


Mathematik: Statistik: Druckversion

<< [Mathematik: Statistik](#)

Druckversion des Buches [Mathematik: Statistik](#).

- Dieses Buch umfasst derzeit etwa 121 DIN-A4-Seiten einschließlich Bilder (Stand: September 2007).
- Wenn Sie dieses Buch drucken oder die Druckvorschau Ihres Browsers verwenden, ist diese Notiz nicht sichtbar.
- Zum Drucken klicken Sie in der linken Menüleiste im Abschnitt „Werkzeuge“ auf [Druckversion](#) .
- Mehr Informationen über Druckversionen siehe [Wikibooks:Druckausgaben](#).
- Hinweise:
 - Für einen reinen Text-Ausdruck kann man die Bilder-Darstellung im Browser deaktivieren.
 - Texte, die in Klappboxen stehen, werden nicht mit ausgedruckt.
 - Die Funktion „Als PDF herunterladen“ kann zu Darstellungsfehlern führen.

Internet-Explorer: Extras > Internetoptionen > Erweitert > Bilder anzeigen (Häkchen entfernen und mit OK bestätigen),
Mozilla Firefox: Extras > Einstellungen > Inhalt > Grafiken laden (Häkchen entfernen und mit OK bestätigen),
Opera: Ansicht > Bilder > Keine Bilder.

Inhalt

- 0. Einführung
- 1. Wahrscheinlichkeitsrechnung
 - 1.1. Zufallsvorgang und Wahrscheinlichkeit
 - 1.2. Gemeinsame Wahrscheinlichkeit mehrerer Ereignisse
 - 1.2.1. Stochastische Unabhängigkeit
 - 1.2.2. Bedingte Wahrscheinlichkeiten
 - 1.2.3. Bayessches Theorem
 - 1.3. Kombinierte Zufallsvorgänge
- 2. Zufallsvariablen
 - 2.1. Diskrete Zufallsvariablen
 - 2.2. Stetige Zufallsvariablen
 - 2.3. Ungleichung von Bienaymé-Tschebyschew
 - 2.4. Mehrdimensionale Zufallsvariablen
 - 2.4.1. Abhängigkeit von Zufallsvariablen
 - 2.4.2. Funktionen von Zufallsvariablen
- 3. Ausgewählte Verteilungen
 - 3.1. Diskrete Verteilungen
 - 3.1.1. Binomialverteilung
 - 3.1.2. Hypergeometrische Verteilung
 - 3.1.3. Poissonverteilung

- 3.2. Stetige Verteilungen
 - 3.2.1. Normalverteilung
 - 3.2.1.1. Zentraler Grenzwertsatz
 - 3.2.1.2. Lineare Funktionen der Normalverteilung
 - 3.2.1.3. Nichtlineare Funktionen der Normalverteilung
 - 3.2.1.3.1. χ^2 -Verteilung
 - 3.2.1.3.2. F-Verteilung
 - 3.2.1.3.3. t-Verteilung
 - 3.3. Approximation von Verteilungen
- 4. Deskriptive Statistik: Analyse eines Merkmals
 - 4.1. Analyse von Daten eines metrischen Merkmals mit wenigen verschiedenen Beobachtungen
 - 4.1.1. Häufigkeitsverteilung
 - 4.1.2. Lageparameter
 - 4.1.3. Streuungsparameter
 - 4.2. Daten, die in Häufigkeitstabellen zusammengefasst sind
 - 4.2.1. Merkmale mit wenig Ausprägungen
 - • • • 4.2.1.1. Verteilung
 - 4.2.1.1.1. Häufigkeit
 - 4.2.1.1.2. Summenfunktion
 - 4.2.1.2. Lageparameter
 - 4.2.1.3. Streuungsparameter
 - 4.2.2. Metrisches Merkmal mit vielen Ausprägungen
 - 4.2.2.1. Verteilung
 - 4.2.2.1.1. Klassierung
 - 4.2.2.1.2. Summenkurve
 - 4.2.2.2. Lageparameter
 - 4.2.2.3. Streuungsparameter
- 5. Mehrere Merkmale
 - 5.1. Korrelationsanalyse
 - 5.2. Regressionsanalyse
 - 5.3. Zeitreihenanalyse
 - 5.3.1. Trend und Saisonkomponente
 - 5.3.2. Glättungsverfahren
- 6. Maßzahlen
 - 6.1. Konzentration
- 7. Schätzen und Testen
 - 7.1. Konfidenzintervalle
 - 7.1.1. Prinzip des Konfidenzintervalls
 - 7.1.2. Ausgewählte Konfidenzintervalle
 - 7.2. Hypothesentests
 - 7.2.1. Verteilungsgebundene Tests
 - 7.2.1.1. Test auf Erwartungswert
 - 7.2.1.2. Test auf Anteilswert
 - 7.2.1.3. Test auf Varianz

- 7.2.2. Verteilungsfreie Tests
 - 7.2.2.1. Prüfung des Zusammenhangs zweier Merkmale
- 8. Übungsaufgaben
- 9. Statistik auf dem Computer
- 10. Literatur
- 11. Tabellen
 - 11.1. Verteilungstabellen
 - 11.1.1. Normalverteilung
 - 11.1.2. Chi-Quadrat-Verteilung
 - 11.1.3. F-Verteilung
 - 11.1.4. t-Verteilung
 - 11.2. Weltdaten
- 12. Sandkasten für Übungsaufgaben

0.	Was ist Statistik?	<< hoch zum Anfang vom Inhaltsverzeichnis << hoch zum Kapitel 0 vom Inhaltsverzeichnis
----	--------------------	---

Was ist Statistik?

Statistik begegnet uns überall im täglichen Leben:

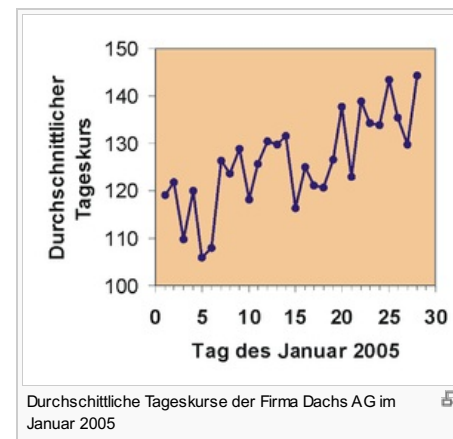
- Die Lebenshaltungskosten sind gegenüber dem Vorjahr um 2 Prozentpunkte gestiegen.
- Im Januar 2005 erzielte die Firma Dachs im Durchschnitt die täglichen Aktienkurse, wie in der Grafik angegeben.
- Hochrechnung von Wahlergebnissen
- Wieviel Gewinn kann eine Lottogesellschaft auswerfen, damit ihr noch Überschuss bleibt?

Was haben diese Beispiele gemeinsam? Sie basieren auf Daten, und zwar sehr vielen Daten. In diese Daten wird Ordnung gebracht: Mit einer Grafik, mit Wahrscheinlichkeiten, mit Durchschnittsberechnungen, mit Vergleichen. Das ist angewandte Statistik.

Wir kommen damit zu einer Definition der Statistik, die relativ kurz und schnörkellos ist, aber im Wesentlichen alles sagt:

Statistik ist die Gesamtheit der Methoden, die für die Untersuchung von Massendaten angewendet werden können.

Ziel der Statistik ist es also, Massendaten zu reduzieren und zu komprimieren, um Gesetzmäßigkeiten und Strukturen in den Daten sichtbar zu machen.



Anwendung im wirtschaftlichen Kontext

Die Lage der Unternehmen heute ist geprägt von Globalisierung, Konkurrenz und Kostendruck. Einsame Manager-Entscheidungen aus dem Bauch heraus führen häufig zum Ruin des Unternehmens. Die Analyse von Wirtschafts- und Unternehmensdaten erlaubt rationale und fundierte Unternehmensentscheidungen. In der Realität sind jedoch Informationen über Unternehmensprozesse nur teilweise bekannt. Gründe dafür sind beispielsweise

1. Die Informationen sind zu komplex, um vollständig erhoben zu werden. **Beispiel:** Der Papierverbrauch in einem großen Unternehmen hängt von vielen Faktoren ab, wie der Zahl der Kopien eines Schreibens, der Neigung der Mitarbeiter, sich alles ausdrücken zu lassen (E-Mails!), dem Umfang des Verteilers für bestimmte Schreiben etc. Man kann den Verbrauch nicht analytisch bestimmen.
2. Zukünftige Unternehmenszahlen sind nicht bekannt und müssen geschätzt werden, z. B. der Cash-Flow einer geplanten Investition für die

Finanzierungsrechnung.

3. Umwelteinflüsse können nicht vorherbestimmt werden, etwa die „Gefahr“ einer Steuererhöhung oder die Akzeptanz eines neuen Produkts durch den Kunden.

In solchen Fällen können keine exakten Entscheidungsgrundlagen geliefert werden. Die resultierenden Lösungen sind „unscharf“. Ein Hilfsmittel für die Entscheidung unter Unsicherheit ist die Statistik.

Da in der heutigen informationsbasierten Gesellschaft eher zu viel als zu wenig Daten verfügbar sind, gewinnt die Statistik als Werkzeug der Entscheidungsfindung immer mehr an Bedeutung.

Einteilung der statistischen Methoden

1. **Deskriptive (beschreibende, empirische) Statistik:** Man untersucht ein Phänomen und fasst die Daten zusammen, ordnet sie, stellt sie grafisch dar. Auf wissenschaftliche Aussagen wird verzichtet.
2. **Induktive (schließende, folgernde, mathematische, analytische) Statistik:** Grundlage ist die Wahrscheinlichkeitstheorie. Ergebnisse der deskriptiven Statistik dienen häufig als Ausgangspunkt für verallgemeinernde Aussagen.

Die mathematische Statistik selbst ist wie die Wahrscheinlichkeitstheorie ein Teilgebiet der Stochastik.

1.	Wahrscheinlichkeitsrechnung	<< hoch zum Anfang vom Inhaltsverzeichnis << hoch zum Kapitel 1 vom Inhaltsverzeichnis
----	-----------------------------	---

Was ist Wahrscheinlichkeit?

Das weiß niemand. Sie ist ein Produkt menschlicher Bemühungen, Ereignisse in der Zukunft vorherzusagen. Sie soll eine Vorstellung über den Grad der Sicherheit vermitteln, mit der ein Ereignis auftritt. Jeder weiß, was es bedeutet, wenn ich sage: Die Wahrscheinlichkeit, eine Sechsz zu würfeln ist größer als die Wahrscheinlichkeit, beim Skat einen Grand zu gewinnen. Aber trotzdem kann man Wahrscheinlichkeit nicht exakt definieren. So könnte man Wahrscheinlichkeitstheorie als Stochern im Nebel bezeichnen. Das hat aber nichts mit dem Begriff Stochastik zu tun!

1.1.	Zufallsvorgang und Wahrscheinlichkeit	<< hoch zum Anfang vom Inhaltsverzeichnis << hoch zum Kapitel 1 vom Inhaltsverzeichnis
------	---------------------------------------	---

Pizzaecken-Beispiel zum Begriff der Wahrscheinlichkeit

Harry und Paula gehen in die Pizzeria. Sie sind frisch verliebt. Paula bestellt sich eine Pizzecke mit Salami und Harry eine mit Schinken. Dann tauschen sie jeweils eine Hälfte, wobei anzumerken ist, dass die Ecken sich in Rand- und Mittelstück teilen lassen. Obwohl Harry normalerweise Randstücke lieber mag, achtet er in seinem aktuellen Zustand nicht darauf. Und auch Paula gibt ihre Hälfte rein nach Zufall ab.

Wie groß ist eigentlich die Wahrscheinlichkeit, dass Harry zwei Randstücke auf dem Teller hat?

Die Meisten antworten richtig: $1/4$.

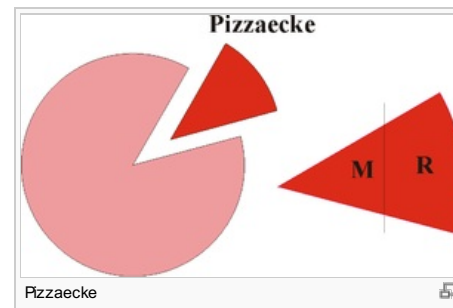
Aber wieso beträgt die Wahrscheinlichkeit ausgerechnet $1/4$?

Betrachten wir den Vorgang:

Bei **gleicher Ausgangslage (Bedingungskomplex)** kann der Versuch, zwei halbe Pizzecken zufällig auszutauschen, **beliebig oft wiederholt** werden. Jeder Versuch hat einen **unsicheren Ausgang**. Es handelt sich hier um einen **Zufallsvorgang (Experiment, Versuch)**.

Der Zufallsvorgang wird also beschrieben durch:

- Gleicher Bedingungskomplex
- Unsicherer Ausgang
- Beliebig oft wiederholbar



Ein bestimmtes Paar Eckhälften auf Harrys Teller ist ein **Ergebnis**. Ein Ergebnis wäre beispielsweise: Die erste Hälfte ist ein Randstück, die zweite Hälfte ist ein Mittelstück,

(R;M) oder kurz RM,

wobei das „linke“ Stück von Harry stammt und das „rechte“ von Paula.

Alle möglichen Paare fasst man in der **Ergebnismenge** Ω zusammen:

$\Omega = \{RR, RM, MR, MM\}$.

Ω ist also die Menge aller möglichen Ergebnisse, die bei einem Zufallsvorgang auftreten können. Führt man diesen Zufallsvorgang unendlich oft durch, müssten vermutlich in 25% aller Versuche zwei Randstücke resultieren, denn man könnte davon ausgehen, dass jedes Paar die gleiche Wahrscheinlichkeit hat, gezogen zu werden. Die Zahl der Ergebnisse, $|\Omega|$ genannt, ist also vier. Deshalb ist die Wahrscheinlichkeit für ein Paar Randstücke

$$P(RR) = \frac{1}{4}.$$

Wenn nun bei einem Versuch beispielsweise „RM“ resultiert, ist das ein **Ereignis**.

Bei „RM“ handelt es sich um ein **Elementarereignis**. Es ist ein Ereignis, das nur **ein** Element der Ergebnismenge enthält.

Es gibt auch kompliziertere, **zusammengesetzte** Ereignisse:

A: Mindestens ein Mittelstück: $A = \{RM, MR, MM\}$

B: Eine komplette Pizzecke: $B = \{RM, MR\}$

Diese Ereignisse beinhalten mehrere Ergebnisse von Ω ; ein Ereignis ist immer eine Teilmenge von Ω .

Das zusammengesetzte Ereignis A tritt also genau dann ein, wenn eines der Elementarereignisse $\{RM\}$, $\{MR\}$ oder $\{MM\}$, die in A enthalten sind, eintritt.

Die Wahrscheinlichkeit als theoretisches Konzept

Kurzer geschichtlicher Überblick

Es werden vermutlich schon so lange Wahrscheinlichkeiten angewendet, wie es den Homo Sapiens gibt. Am letzten Tag der Schlacht im Teutoburger Wald (9 n. Chr.) gab es ein Gewitter. Die Römer deuteten es als warnenden Hinweis von Merkur, des Gottes von Blitz und Donner. Die Germanen sahen es als Aufmunterung des Kriegsgottes Thor. Wie man weiß, hatten beide Parteien recht.

Im 17. Jahrhundert, dem Zeitalter des Rationalismus, befasste sich **Blaise Pascal** (1623 - 1662) systematisch mit Wahrscheinlichkeiten im Glücksspiel und begründete so die Wahrscheinlichkeitsrechnung als eigenständige Disziplin.

Jakob Bernoulli (1654 - 1705) befasste sich ebenfalls mit Fragen der diskreten Wahrscheinlichkeiten und gab vermutlich das erste Buch über Wahrscheinlichkeitsrechnung heraus.

Mit **Abraham de Moivre** (1667 - 1754) und **Pierre Simon Laplace** (1749 - 1827) wurde bereits die Normalverteilung entwickelt und von **Carl Friedrich Gauß** (1777 – 1855) weiter bearbeitet.

Richard Edler von Mises (1883 - 1953) lieferte wertvolle Beiträge zur Schätzung von Wahrscheinlichkeiten und zur mathematischen Statistik.

1933 schlug der russische Mathematiker **Andrej Nikolajewitsch Kolmogorow** (1903 - 1987) eine **axiomatische Definition der Wahrscheinlichkeit** vor, auf der die heutige Wahrscheinlichkeitstheorie basiert. Diese Definition ist eine Anwendung der Maßtheorie.

Ergebnisse und Ereignisse

Das heutige Konzept der Wahrscheinlichkeitsrechnung präsentiert sich folgendermaßen:

Gegeben ist die **Ergebnismenge** (Ereignisraum, Stichprobenraum) Ω eines **Zufallsvorgangs**. Diese Menge enthält alle möglichen Ergebnisse, die ein Zufallsvorgang hervorbringen kann. Je nach Art des Zufallsvorgangs muss man verschiedene Ergebnismengen betrachten:

Ω enthält **endlich** viele Ergebnisse.

Beispiele:

- Zufallsvorgang: 1x Würfeln. $\Omega = \{1, 2, 3, 4, 5, 6\}$.

- Zufallsvorgang: Augenfarbe der nächsten Person, die bei einem Casting vorspricht. $\Omega = \{\text{blau, grün, braun}\}$.

Ω enthält **abzählbar unendlich** viele Ergebnisse.

Beispiele:

- Zufallsvorgang: Zahl der Autos, die eine Stunde lang ab 12 Uhr bei einer Fahrzeugzählung an einer bestimmten Zählstelle vorbeifahren. $\Omega = \{0, 1, 2, 3, \dots\}$.
- Zufallsvorgang: Zahl der Anforderungen an einen Server innerhalb einer Stunde. $\Omega = \{0, 1, 2, \dots\}$.

Man kann zwar die Ergebnisse durchzählen, aber es kann keine vernünftige Obergrenze angegeben werden, deshalb lässt man die Obergrenze offen.

Ist Ω weder abzählbar noch abzählbar unendlich, so enthält Ω **überabzählbar** viele Ergebnisse. Man könnte auch sagen, die Ergebnismenge ist ein Intervall der reellen Zahlen.

Beispiele:

- Zufallsvorgang: Eine erwachsene Person wird gewogen (in kg). $\Omega = \{x | 30 \leq x \leq 200; x \in \mathbb{R}\}$.
- Zufallsvorgang: Cash-Flow eines Unternehmens (in €). $\Omega = \mathbb{R}$.

Cash-Flow bezeichnet übrigens die Differenz Einnahmen - Ausgaben, bzw. präziser: Einzahlungen - Auszahlungen.

Hier können die Ergebnisse nicht mehr abgezählt werden. Ein beliebig kleines Intervall der Ergebnismenge enthält unendlich viele Elemente. Was ist das nächstgrößere Element von 50 kg: 51 kg, 50,01 kg oder 50,000000001 kg? Im Intervall [50, 51] sind also unendlich viele Elemente.

Man könnte hier einwenden, dass doch beispielsweise Cash-Flow als kleinste Einheit Cent hat, also doch eigentlich abzählbar ist. Das stimmt natürlich, aber bei sehr vielen, nah zusammenliegenden Elementen vereinfacht man die Analyse, indem man die Menge als stetig annimmt. Man spricht hier von Quasistetigkeit.

Hat ein Zufallsvorgang ein konkretes Ergebnis erbracht, ist ein **Ereignis** eingetreten. Es gibt einfache Ereignisse, die lediglich ein Ergebnis enthalten, so genannte **Elementarereignisse** und es gibt komplexere Ereignisse, die sich aus mehreren Ergebnissen zusammensetzen. **Ein Ereignis A ist immer eine Teilmenge der Ergebnismenge Ω .**

Da Ereignisse Mengen sind, können alle Operationen der **Mengenalgebra**, die mit der **Booleschen Algebra** (auch Schaltalgebra) gleichgesetzt werden kann, angewendet werden. Grundlegende Operationen für Mengen der Booleschen Algebra sind " („nicht“ als Komplement), \cap und \cup . Alle anderen Operationen können daraus hergeleitet werden.

Alle interessierenden Ereignisse fasst man nun in einer so genannten **Ereignismenge (Ereignissystem) E** zusammen. **E** ist also eine Menge von Teilmengen. Damit diese Menge mit der Booleschen Algebra bearbeitet werden kann, muss sie entsprechende Forderungen erfüllen:

- Wenn das Ereignis A in E enthalten ist, muss auch sein Komplement \bar{A} enthalten sein.
- Wenn A und B enthalten sind, muss auch $A \cup B$ enthalten sein (Man kann ausrechnen, dass dann auch $A \cap B$ enthalten ist).
- Es muss das „Null-Element“ \emptyset enthalten sein (Das impliziert, dass auch „1-Element“ Ω , welches das Komplement von \emptyset ist, enthalten ist).

Die umfassendste Ereignismenge ist die Potenzmenge P , die alle Teilmengen von Ω enthält.

Beispiel einer Potenzmenge:

Zufallsvorgang: Aus einer Urne mit einer blauen (b), einer roten (r) und einer gelben (g) Kugel wird eine Kugel gezogen. Wir interessieren uns für die Farbe der Kugel.

Ergebnismenge: $\Omega = \{g, b, r\}$

Potenzmenge: $P = \{\emptyset, \{r\}, \{g\}, \{b\}, \{r, g\}, \{r, b\}, \{g, b\}, \{r, g, b\}\}$

Ausgehend von dieser Konstellation hat **Kolmogorow** mit seinen Axiomen ein **Wahrscheinlichkeitsmaß** konstruiert, d.h. eine Abbildung der Ergebnismenge Ω auf die Menge der reellen Zahlen im Intervall [0;1]:

$$F: \Omega \rightarrow \mathbb{R}; A \rightarrow P(A)$$

Eine Funktion P , die jedem Ereignis A aus \mathcal{E} eine reelle Zahl zuordnet, heißt Wahrscheinlichkeit, wenn sie folgende Axiome erfüllt:

Axiome der Wahrscheinlichkeiten:

Gegeben sind zwei Ereignisse $A, B \subset \Omega$.

1. $P(A) \geq 0$. **Nichtnegativität**
2. $P(\Omega) = 1$. **Normiertheit**
3. $P(A \cup B) = P(A) + P(B)$, falls A und B disjunkt sind. **Additivität**

Dieses Axiomensystem kann nur auf endlich viele Ereignisse angewendet werden. Für unendlich viele Ereignisse A_i ($i = 1, 2, \dots$) erhält man statt der endlichen Ereignismenge die **σ -Algebra**. Sie enthält alle geforderten Eigenschaften der Ereignismenge auf unendlich viele Ereignisse A_i ausgeweitet. Hier wird das 3. Axiom entsprechend angepasst:

3. Sind die Ereignisse A_i sämtlich paarweise disjunkt, ist bei ihrer Vereinigung

$$P(A_1 \cup A_2 \cup A_3 \dots) = P(A_1) + P(A_2) + P(A_3) + \dots, (\sigma\text{-Additivität}).$$

Berechnung der Wahrscheinlichkeit eines Ereignisses

Es müssen nun noch die Ereignisse mit Wahrscheinlichkeiten ausgestattet werden. Auf welche Weise das geschehen soll, ist in den Axiomen nicht angegeben. Es gibt hier verschiedene Verfahren. Man erhält schließlich die Wahrscheinlichkeitsverteilung.

Wie ordnen wir den Ereignissen am besten Wahrscheinlichkeiten zu?

Betrachten wir im Pizzatecken-Beispiel das Ereignis A : Mindestens ein Mittelstück. Es ist $A = \{RM, MR, MM\}$. A belegt in Ω drei von vier möglichen Ergebnissen, also ist die Wahrscheinlichkeit $P(A) = 3/4$. Diese Vorgehensweise entspricht der **Klassischen**

Wahrscheinlichkeitsauffassung. Man bezeichnet sie als **Symmetrieprinzip** oder **Prinzip nach LAPLACE**:

Jedes Ergebnis ist gleich häufig. $|A|$ ist die Zahl der Ergebnisse, die durch A belegt werden (Anzahl der günstigen Ergebnisse), $|\Omega|$ ist die Zahl aller möglichen Ergebnisse. Es ist

$$P(A) = \frac{|A|}{|\Omega|} = \frac{3}{4}.$$

Das Symmetrieprinzip hat allerdings den Nachteil, dass es nicht bei allen Zufallsvorgängen angewendet werden kann, z.B. bei unendlich vielen Ergebnissen. Oft ordnet man auch Ergebnissen unterschiedliche Wahrscheinlichkeiten zu, z.B.

Zufallsvorgang: Wetter von heute.

Ergebnismenge $\Omega = \{\text{schön, schlecht}\}$.

$P(\text{„schön“}) = 0,6$, $P(\text{„schlecht“}) = 0,4$.

Wie kommt man auf diese Wahrscheinlichkeiten 0,4 und 0,6? Man hat in diesem Fall etwa die Wetteraufzeichnungen der letzten 100 Jahre ausgewertet und hat festgestellt, dass der Anteil der schönen Tage 60 % betrug. Wir haben hier eine Anwendung der **Statistischen**

Wahrscheinlichkeitsauffassung: Man führt ein Zufallsexperiment sehr oft durch. Mit steigender Zahl der Versuche nähert sich der Anteil der Versuche, die das Ereignis A hervorgebracht haben, der „wahren“ Wahrscheinlichkeit $P(A)$, formal ausgedrückt

$$P(A) = \lim_{n \rightarrow \infty} \frac{n(A)}{n},$$

mit $n(A)$ als Zahl der Versuche, die das Ereignis A hervorgebracht haben. Man bezeichnet diesen Zusammenhang als **Gesetz der großen Zahlen**. Er liefert die Begründung, dass man unbekannte Wahrscheinlichkeiten mit Hilfe von empirischen Beobachtungen schätzen kann, wobei hier gilt: Viel hilft viel!

Bei manchen Fragestellungen versagen die beiden obigen Wahrscheinlichkeitskonzepte. Z.B. bei Ereignissen, die sehr selten auftreten, für die

man also auch keine Versuchsreihen zur Verfügung hat, etwa die Wahrscheinlichkeit für den Erfolg eines neu auf dem Markt platzierten Produkts. Es möchte beispielsweise ein Unternehmen ein neues Spülmittel auf den Markt bringen. Es steht vor der Alternative, Fernsehwerbung einzusetzen oder nicht. Es ist mit den Ereignissen konfrontiert: Wenn Fernsehwerbung eingesetzt wird, ist das Spülmittel ein Erfolg/kein Erfolg. Wenn keine Fernsehwerbung eingesetzt wird, ist das Spülmittel ein Erfolg/kein Erfolg. Für diese vier Ereignisse sollen Wahrscheinlichkeiten ermittelt werden. Da man keine verlässlichen Informationen darüber hat, wird man aus dem Bauch heraus, eventuell unter Berücksichtigung ähnlicher Erfahrungen bestimmte Wahrscheinlichkeiten zuordnen. Dieses Vorgehen entspricht der **Subjektiven Wahrscheinlichkeitsauffassung**.

Da Ereignisse als Mengen definiert sind, kann man auch in vielen Fällen Ereignisse und ihre Wahrscheinlichkeiten in **Venn-Diagrammen** veranschaulichen. Die Wahrscheinlichkeit ist dann die Fläche der entsprechenden Menge. Manchmal ist es hilfreich, das Venn-Diagramm maßstabsgetreu auf kariertes Papier abzutragen, indem die Mengen rechteckig dargestellt werden.

Pizzeria-Beispiel zur Berechnung von Wahrscheinlichkeiten

Jetzt schauen wir uns in der Pizzeria etwas genauer um: Der Inhaber Carlo Pomodoro ist ein mitleidiger Mensch und duldet auch arme Gäste, die sich nichts bestellen. Deshalb ist das Lokal mit seinen 50 Gästen eigentlich schon überfüllt. 20 Personen haben sich Pizza bestellt und 10 Lasagne. Das Essen ist so reichlich, dass niemand zwei Mahlzeiten bestellt. 40 Gäste trinken Wein und 20 Gäste trinken Mineralwasser, aber 15 trinken Wasser und Wein.

Wir ziehen zufällig einen Gast aus der fröhlich lärmenden Menge. Wie groß ist die Wahrscheinlichkeit, einen Pizza-Esser zu erhalten?

Wir haben $|\Omega| = 50$ verschiedene Ergebnisse. Man kann davon ausgehen, dass jeder Gast die gleiche Wahrscheinlichkeit hat, gezogen zu werden.

Wir definieren nun die Ereignisse:

- A: Der Gast isst Pizza; B: Der Gast isst Lasagne;
- C: Der Gast trinkt Wein; D: Der Gast trinkt Wasser.

Nach dem Symmetrieprinzip ist

$$P(A) = \frac{|A|}{|\Omega|} = \frac{20}{50} = \frac{2}{5},$$

$$P(B) = \frac{10}{50} = \frac{1}{5},$$

$$P(C) = \frac{4}{5} \text{ und } P(D) = \frac{2}{5}.$$

Wir können berechnen:

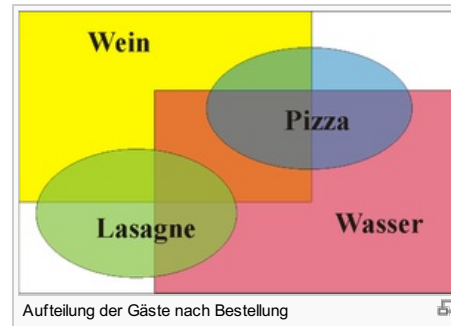
Wahrscheinlichkeit, dass jemand Wasser **und** Wein trinkt:

$$P(C \cap D) = \frac{|C \cap D|}{|\Omega|} = \frac{15}{50} = \frac{3}{10}.$$

Wahrscheinlichkeit, dass ein zufällig ausgewählter Gast kein Wasser trinkt (\bar{D}):

$$P(\bar{D}) = \frac{|\bar{D}|}{|\Omega|} = \frac{50 - 20}{50} = 1 - \frac{20}{50} = \frac{3}{5} = 1 - P(D).$$

Anteil der Leute, die Wasser **oder** Wein trinken:



$$P(C \cup D) = P(C) + P(D) - P(C \cap D) = \frac{40}{50} + \frac{20}{50} - \frac{15}{50} = \frac{45}{50} = \frac{9}{10}.$$

Diese Beziehung gilt immer für zwei Ereignisse!

Wahrscheinlichkeit, dass ein Gast Pizza oder Lasagne isst:

$$P(A \cup B) = P(A) + P(B) - P(A \cap B) = \frac{20}{50} + \frac{10}{50} - 0 = \frac{30}{50}.$$

Die Mengen A und B sind disjunkt.

Wahrscheinlichkeit, dass der zufällig ausgewählte Gast kein Wasser oder keinen Wein trinkt:

$$P(\bar{C} \cup \bar{D}) = P(\bar{C}) + P(\bar{D}) - P(\bar{C} \cap \bar{D}).$$

Hier ist die direkte Berechnung der Wahrscheinlichkeit analog zu oben umständlich. Man verwendet am besten die

DE MORGANSche Regel:

$$P(\bar{C} \cup \bar{D}) = P(\overline{C \cap D}) = 1 - P(C \cap D) = 1 - \frac{15}{50} = \frac{35}{50} = 0,7.$$

Was gelernt werden muss

Ein Ereignis $A (A \subset \Omega)$:

$$0 \leq P(A) \leq 1.$$

$$P(\bar{A}) = 1 - P(A).$$

$$P(\emptyset) = 0.$$

Zwei Ereignisse A und B ($A, B \subset \Omega$) :

A und B sind im allgemeinen nicht disjunkt, also ist die Wahrscheinlichkeit, dass A **oder** B eintritt, nach dem **Additionssatz für zwei Ereignisse**:

$$P(A \cup B) = P(A) + P(B) - P(A \cap B).$$

Falls A und B **disjunkt** sind, ist

$$P(A \cup B) = P(A) + P(B).$$

DE MORGANSche Regeln:

$$P(\bar{A} \cup \bar{B}) = P(\overline{A \cap B})$$

und

$$P(\bar{A} \cap \bar{B}) = P(\overline{A \cup B})$$

Für drei Ereignisse A_i ($i=1, 2, 3$) aus Ω gilt analog zu obigen Überlegungen:

$$P(A_1 \cup A_2 \cup A_3) = P(A_1) + P(A_2) + P(A_3) - P(A_1 \cap A_2) - P(A_1 \cap A_3) - P(A_2 \cap A_3) + P(A_1 \cap A_2 \cap A_3).$$

Mehrere Ereignisse A_i (i endlich oder unendlich):

Sind die Ereignisse A_i sämtlich paarweise disjunkt, ist bei ihrer Vereinigung

$$P(A_1 \cup A_2 \cup A_3 \dots) = P(A_1) + P(A_2) + P(A_3) + \dots$$

Übung

Zeigen Sie anhand eines Venn-Diagramms die Gültigkeit einer der DeMorganschen Regeln.

1.2.

Gemeinsame Wahrscheinlichkeit mehrerer Ereignisse

<< hoch zum Anfang vom Inhaltsverzeichnis

<< hoch zum Kapitel 1 vom Inhaltsverzeichnis

Stochastische Unabhängigkeit

Ein häufiges Untersuchungsobjekt in der Statistik ist, ob verschiedene Ereignisse **abhängig** oder **unabhängig** voneinander sind, d.h. ob das Zustandekommen eines Ereignisses durch ein anderes begünstigt wird. So untersucht man beispielsweise in der Marktforschung, ob Status und Bildung eines Konsumenten die Ausgaben für eine bestimmte Zeitschrift beeinflussen.

Beispiel zum Begriff der stochastischen Unabhängigkeit

Eine umfangreiche Marketingstudie über Zahnputzgewohnheiten von Konsumenten hat ergeben, dass 50 % der Studierenden einer kleinen Hochschule bei ihren Eltern wohnen. Ebenso, dass 50 % der Studierenden Zahnpasta mit roten Streifen und 50 % andersfarbige Zahnpasta bevorzugen.

Betrachten wir den **Zufallsvorgang**: Eine Studentin kommt in einen Laden und kauft Zahnpasta. Es seien definiert die Ereignisse:

E: Die Studentin wohnt bei ihren Eltern.

R: Die Studentin kauft Zahnpasta mit roten Streifen.

Frage: Hat der Wohnort der Studentin einen Einfluss auf die Farbpräferenz?

Vermutlich nein, die Ereignisse E und R sind **stochastisch unabhängig**, d.h. in wahrscheinlichkeitstheoretischer Hinsicht unabhängig.

Wir interessieren uns zunächst für den Wohnort der Studierenden. In der Grafik 1 ist die Ergebnismenge nach dem Wohnort aufgeteilt.

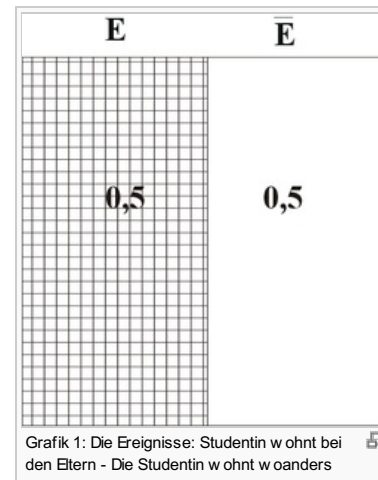
Frage: Wieviel Prozent der Studierenden wohnen bei ihren Eltern **und** werden voraussichtlich Zahnpasta mit roten Streifen kaufen?

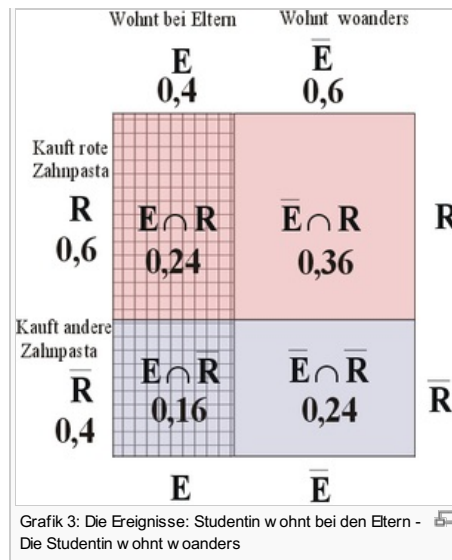
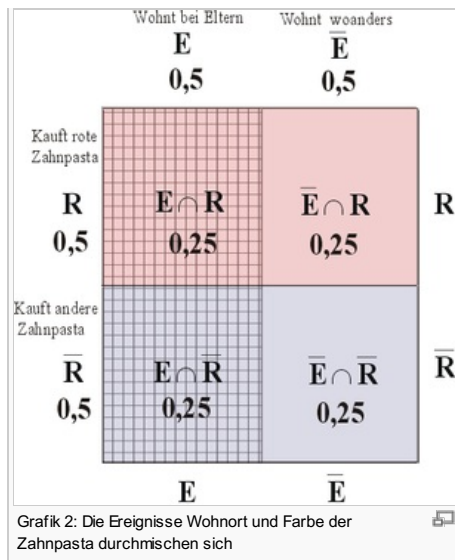
Da sich bei Unabhängigkeit der Ereignisse die Studierenden in Bezug auf ihre Farbpräferenz gleichmäßig auf die Wohnorte verteilen, werden wohl 50 % der Rotkäufer bei ihren Eltern wohnen und 50 % woanders. D.h. 50 % von 50 % der Studierenden wohnen bei ihren Eltern **und** bevorzugen rote Zahnpasta. Es gilt also:

$$P(R \cap E) = 0,5 \cdot 0,5 = 0,25.$$

Die Grafik 2 zeigt, wie sich bei Unabhängigkeit der Variablen Wohnort und Farbpräferenz die Wahrscheinlichkeiten der Farbpräferenz auf die Wohnorte aufteilen.

Ist nun beispielsweise $P(E) = 40\%$ und $P(R) = 60\%$, ergibt sich bei Unabhängigkeit die Aufteilung wie in der Grafik 3, denn auch hier müssten 60 % der „Nesthocker“ und 60 % der „Nestflüchter“ gleichermaßen Zahnpasta mit roten Streifen kaufen.





Beispiel zum Begriff der stochastischen Abhängigkeit

Oben haben wir den Fall betrachtet, dass zwei Ereignisse unabhängig sind. Im Allgemeinen muss man aber davon ausgehen, dass Ereignisse, die man gemeinsam analysiert, abhängig sind.

Im Rahmen der Marketingstudie wurden Daten eines Gesundheitsamtes in Musterstadt verwendet, die die Zahngesundheit von Schulkindern betraf. Man weiß aus dieser Studie, dass 50 % der Schulkinder Karies haben und 50 % der Schulkinder sich regelmäßig die Zähne putzen.

Wir betrachten den Zufallsvorgang: Es wird ein Schulkind zufällig ausgewählt.

Wir definieren als Ereignisse

Z: Das Schulkind putzt sich regelmäßig die Zähne.

K: Das Schulkind hat Karies.

Ist nun

$$P(Z \cap K) > P(Z \cap \bar{K})$$

oder $P(Z \cap K) < P(Z \cap \bar{K})$?

Ist also die Wahrscheinlichkeit, ein Kind zu erhalten, das sich regelmäßig die Zähne putzt **und** Karies hat, größer als die Wahrscheinlichkeit, ein Kind zu erhalten, das sich regelmäßig die Zähne putzt **und** keine Karies hat, oder ist es umgekehrt, oder sind vielleicht die Wahrscheinlichkeiten gleich?

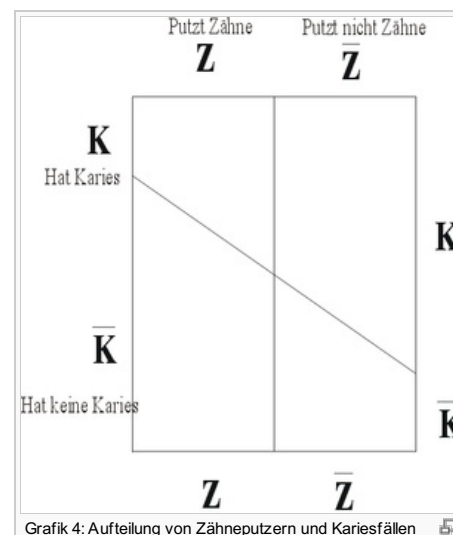
Es ist vermutlich

$$P(Z \cap K) < P(Z \cap \bar{K}),$$

denn Zähneputzen und Karies sind bekanntlich nicht unabhängig voneinander zu betrachten. Also sind Z und K stochastisch abhängige Ereignisse. Wir werden vermutlich eine Aufteilung der gemeinsamen Wahrscheinlichkeiten erhalten, die ähnlich der Grafik 4 ist. Besonders groß sind $P(Z \cap \bar{K})$ und $P(\bar{Z} \cap K)$.

Die gemeinsamen Wahrscheinlichkeiten können allerdings nicht mit unseren Informationen bestimmt werden, sie hängen von der Stärke der Abhängigkeit ab.

Bei **stochastisch abhängigen Ereignissen** interessiert man sich häufig für das



bedingte Auftreten eines Ereignisses, z.B. für die **bedingte Wahrscheinlichkeit**

$$P(K|\bar{Z}),$$

dass ein zufällig ausgewähltes Schulkind Karies hat, wenn man weiß, dass es sich nicht regelmäßig die Zähne putzt.

Bedingte Wahrscheinlichkeiten

Beispiel

Einige Jahre später wurde in der Grundschule von Musterdorf zu Forschungszwecken wieder an 200 Kindern eine Reihenuntersuchung zur Zahngesundheit durchgeführt. Jetzt putzten sich 60 % der Kinder regelmäßig die Zähne. Von diesen Kindern hatten 40 Karies. Bei den Zahnputzmuffeln hatten 60 Kinder Karies.

Wir wollen ein maßstabsgetreues Venndiagramm konstruieren. Jedes Kästchen steht für 5 Kinder. Es sind

$$\begin{aligned} P(Z) &= 0,6; & P(\bar{Z}) &= 0,4; \\ P(Z \cap K) &= 0,2; & P(Z \cap \bar{K}) &= 0,4; \\ P(\bar{Z} \cap K) &= 0,3; & P(\bar{Z} \cap \bar{K}) &= 0,1. \end{aligned}$$

Wir interessieren uns nun für die bedingte Wahrscheinlichkeit, dass ein Kind Karies hat, wenn bekannt ist, dass es sich die Zähne putzt:

$$P(K|Z).$$

In andere Worte gekleidet: **Der Anteil der Kinder mit Karies an den Kindern, die sich regelmäßig die Zähne putzen.**

Es gilt für die bedingte Wahrscheinlichkeit

$$P(K|Z) = \frac{P(K \cap Z)}{P(Z)}.$$

Wie ist diese Wahrscheinlichkeit zu verstehen?

Es werden zunächst alle Kinder, die sich regelmäßig die Zähne putzen, in die Aula geschickt. Aus diesen 120 Kindern wird nun zufällig eins ausgewählt. Mit welcher Wahrscheinlichkeit hat dieses Kind Karies? Wir betrachten also 120 zahnputzende Kinder, davon haben 40 Kinder Karies.

Genau diese Vorgehensweise ist das Prinzip der bedingten Wahrscheinlichkeiten!

$$\text{Es ergibt sich: } P(K|Z) = \frac{40}{120} = \frac{1}{3}.$$

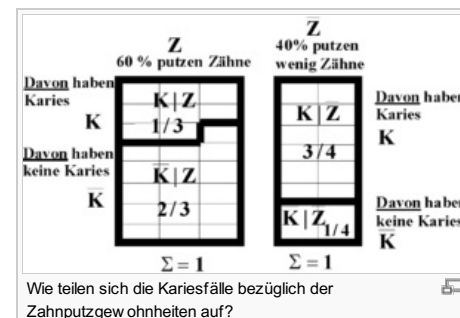
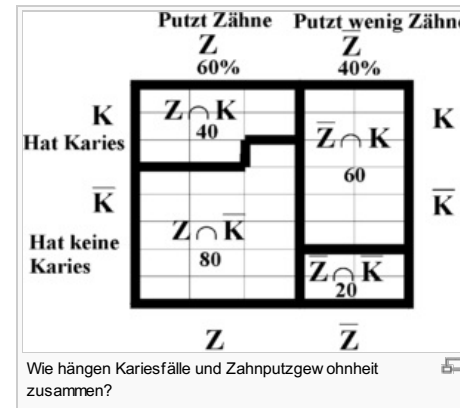
Ein Drittel der zahnputzenden Kinder hat Karies: Dann haben natürlich zwei Drittel der zahnputzenden Kinder keine Karies. Wir sehen sogleich, dass die obige Rechnung die schon bekannte Formel

$$P(K|Z) = \frac{P(K \cap Z)}{P(Z)} = \frac{\frac{40}{200}}{\frac{120}{200}} = \frac{40}{120} = \frac{1}{3},$$

darstellt. Entsprechend erhalten wir

$$P(\bar{K}|Z) = \frac{P(\bar{K} \cap Z)}{P(Z)} = \frac{\frac{80}{200}}{\frac{120}{200}} = \frac{80}{120} = \frac{2}{3},$$

$$P(K|\bar{Z}) = \frac{P(K \cap \bar{Z})}{P(\bar{Z})} = \frac{\frac{60}{200}}{\frac{80}{200}} = \frac{3}{4},$$



$$P(\bar{K}|\bar{Z}) = \frac{P(\bar{K} \cap \bar{Z})}{P(\bar{Z})} = \frac{\frac{20}{200}}{\frac{80}{200}} = \frac{1}{4}.$$

Vergleichen Sie das Venndiagramm mit dem vorhergehenden! Wieso unterscheiden sich beide Diagramme?

Übung

Es ist bekannt, dass die Aktienkurse des Unternehmens Dachs an 55% aller Börsentage gestiegen sind.

Ereignisse: K_1 : Der Kurs steigt am ersten Tag K_2 : Der Kurs steigt am zweiten Tag

Man hat folgende Gesetzmäßigkeit der Kursentwicklung festgestellt: In 40 % aller Beobachtungen stieg der Kurs am ersten Tag und am zweiten Tag, in 15 % der Beobachtungen stieg der Kurs am ersten Tag und fiel am zweiten Tag. Dagegen fiel in 15 % der Beobachtungen der Kurs am ersten Tag und stieg am zweiten Tag. An den restlichen Tagespaaren fiel der Kurs an beiden Tagen.

1. Stellen Sie die gemeinsamen Wahrscheinlichkeiten im Venndiagramm grafisch dar.
2. Sind die Ereignisse K_1 und K_2 stochastisch unabhängig? (Begründen Sie die Antwort formal mit Hilfe der Wahrscheinlichkeitstheorie.)
3. Am heutigen Tag ist der Kurs gestiegen.
 - Mit welcher Wahrscheinlichkeit wird er morgen steigen (Gesucht: $P(K_2|K_1)$)?
 - Mit welcher Wahrscheinlichkeit wird er dagegen fallen?
4. Mit welcher Wahrscheinlichkeit wird der Kurs morgen steigen, wenn er heute gefallen ist?

Bayessches Theorem

Häufig liegen die Informationen über zwei Ereignisse nur als bedingte Wahrscheinlichkeiten vor. Wie kann man sie weiter verwenden?

Beispiel für zwei Ereignisse

Ein bekannter Vergnügungspark verbraucht täglich große Mengen an Glühbirnen für die Dekoration der Stände. Damit die Verbrauchskosten nicht so hoch werden, setzen sich die Glühbirnen nur zu 60% aus Markenware und zu 40 % aus markenfreier Ware zusammen. Aufgrund langjähriger Beobachtungen weiß man, dass von den Marken-Glühbirnen pro Monat 5% defekt werden. Jedoch werden von den markenfreien Glühbirnen monatlich 10% defekt.

Zunächst wollen wir das Gegebene grafisch (Grafik 5) darstellen: Wenn von den Markenglühbirnen 5 % defekt werden, bleiben 95% heil. 5% ist also **Anteil** der defekten Glühbirnen an den Markenglühbirnen, d.h. es handelt sich um die bedingte Wahrscheinlichkeit $P(D|M)$ usw.

Der Betreiber des Vergnügungsparks braucht für die Kostenplanung des nächsten Sommers die Information, wie groß der Anteil der Markenglühbirnen an den defekten Glühbirnen ist, d.h. er sucht $P(M|D)$. Das bedeutet: **Alle defekten Glühbirnen eines Tages werden in einem Korb gesammelt. Es wird eine Glühbirne zufällig entnommen.** Mit welcher Wahrscheinlichkeit erhält man eine Markenbirne?

Wir wissen, dass gilt:

$$P(M|D) = \frac{P(M \cap D)}{P(D)}.$$

Leider sind aber die Komponenten des Bruchs unbekannt. Wir werden nun eine Methode finden, sie doch zu berechnen.

Zunächst suchen wir den Zähler $P(M \cap D)$: Wir kennen $P(D|M)$. Bekanntlicherweise berechnet es sich als

$$P(D|M) = \frac{P(M \cap D)}{P(M)}.$$

Also ist der gesuchte Zähler auch in $P(D|M)$ enthalten und kann ganz einfach durch Auflösung der Gleichung berechnet werden als

$$P(M \cap D) = P(D|M)P(M).$$

also

$$P(M \cap D) = 0,05 \cdot 0,6 = 0,03.$$

	M 60 % Marken	\bar{M} 40% NoName	
Davon defekt D	D M 0,05	D \bar{M} 0,1	Davon defekt D
Davon OK \bar{D}	\bar{D} M 0,95	\bar{D} \bar{M} 0,9	Davon OK \bar{D}
	$\Sigma = 1$	$\Sigma = 1$	

Grafik 5

Jetzt fehlt noch der Nenner $P(D)$. Betrachten wir das Venndiagramm Grafik 6. D setzt sich aus den Schnittmengen $D \cap M$ und $D \cap \bar{M}$ zusammen.

Die gesamte Wahrscheinlichkeit von D ist also die Summe

$$P(D) = P(M \cap D) + P(\bar{M} \cap D).$$

eine Erkenntnis, die man auch als Satz der totalen Wahrscheinlichkeit bezeichnet, und das gibt, wie wir oben gesehen haben,

$$P(D) = P(D|M)P(M) + P(D|\bar{M})P(\bar{M}),$$

in unserem Beispiel

$$P(D) = 0,05 \cdot 0,6 + 0,1 \cdot 0,4 = 0,07.$$

Es sind also 7% aller Glühbirnen defekt.

Die gesuchte bedingte Wahrscheinlichkeit ist nun

$$P(M|D) = \frac{P(M \cap D)}{P(D)} = \frac{P(D|M)P(M)}{P(D|M)P(M) + P(D|\bar{M})P(\bar{M})}.$$

Diese Formel wird als Bayessches Theorem bezeichnet.

Die gesuchte Wahrscheinlichkeit beträgt

$$P(M|D) = \frac{0,03}{0,07} = 0,4286.$$

Diese Wahrscheinlichkeit fällt deshalb so überraschend hoch aus, weil 50% mehr Markenbirnen als markenfreie verwendet werden. Entsprechend ist der Anteil der markenfreien Glühbirnen an den defekten 0,5714.

Wir wollen nun mehr als zwei Ereignisse analysieren.

Beispiel für mehr als zwei Ereignisse

Eine Spedition beschäftigt drei LKW-Fahrer, die Herren Ahorn, Behorn und Zehorn. Ahorn fährt 50% aller Fuhren, Behorn 20% und Zehorn 30%. Aus Erfahrung weiß man, dass Ahorn bei 10% aller Fahrten eine Beule verursacht, Behorn bei 15% aller Fahrten und Zehorn bei 20% aller Fahrten (Grafik 7).

Wir definieren die Ereignisse:

F_1 : Ahorn ist gefahren, F_2 : Behorn ..., F_3 : Zehorn ...

B: Eine Beule wurde gefahren.

Wir wollen zuerst das Gegebene festhalten: Wenn Ahorn in 10 % aller Fahrten eine Beule fährt, wickelt er die restlichen 90 % ohne Schaden ab usw.

Man interessiert sich für die Wahrscheinlichkeit, dass Ahorn gefahren ist, wenn wieder ein Mal eine Beule in einem LKW auftaucht, d.h. für $P(F_1|B)$.

Es ist wieder

$$P(F_1|B) = \frac{P(F_1 \cap B)}{P(B)}.$$

Nach dem Multiplikationssatz der Wahrscheinlichkeiten muss

$$P(F_1 \cap B) = P(B|F_1)P(F_1)$$

sein, also

$$P(F_1 \cap B) = 0,1 \cdot 0,5 = 0,05.$$

Aber wie erhalten wir $P(B)$? Auch hier gilt wieder der Satz von der totalen Wahrscheinlichkeit, z.B.:

	M	\bar{M}	
D	$D \cap M$ 0,03	$D \cap \bar{M}$ 0,04	D $P(D) =$ 0,04 + 0,03
\bar{D}	$\bar{D} \cap M$ 0,57	$\bar{D} \cap \bar{M}$ 0,36	
	$\Sigma = 1$		

Grafik 6

	F_1 50%	F_2 20%	F_3 30%
Dabei Beule B	0,1	0,15	0,2
Dabei keine Beule \bar{B}	0,9	0,85	0,8

Grafik 7

$$P(F_1 \cap B) = P(B|F_1) \cdot P(F_1).$$

Wir erhalten dann für P(B)

$$P(B) = P(F_1 \cap B) + P(F_2 \cap B) + P(F_3 \cap B) \\ = P(B|F_1)P(F_1) + P(B|F_2)P(F_2) + P(B|F_3)P(F_3),$$

also

$$P(B) = 0,1 \cdot 0,5 + 0,15 \cdot 0,2 + 0,2 \cdot 0,3 = 0,05 + 0,03 + 0,06 = 0,14.$$

Unsere gesuchte Wahrscheinlichkeit beträgt

$$P(F_1|B) = \frac{P(F_1 \cap B)}{P(B)} = \frac{0,05}{0,14} = 0,3571.$$

Entsprechend sind

$$P(F_2|B) = \frac{0,03}{0,14} = 0,2143$$

und

$$P(F_3|B) = \frac{0,06}{0,14} = 0,4286$$

Also hat Zehorn mit größter Wahrscheinlichkeit die Beule gefahren.

Wir fassen nun das Gelernte dieser Seite zusammen:

Theoretische Erkenntnisse

Zwei Ereignisse A und B aus Ω :

Sind zwei Ereignisse A und B **stochastisch unabhängig**, ist ihre gemeinsame Wahrscheinlichkeit gleich dem Produkt der Einzelwahrscheinlichkeiten:

$$P(A \cap B) = P(A) \cdot P(B).$$

Man beachte: Ereignisse sind grundsätzlich **nicht** als unabhängig zu betrachten!

Die **bedingten Wahrscheinlichkeiten** für A und B sind

$$P(A|B) = \frac{P(A \cap B)}{P(B)} \text{ und } P(B|A) = \frac{P(A \cap B)}{P(A)}.$$

Allgemeiner Multiplikationssatz der Wahrscheinlichkeiten:

$$P(A \cap B) = P(A|B)P(B) = P(B|A)P(A).$$

Theorem von BAYES:

$$P(B|A) = \frac{P(A \cap B)}{P(A)} = \frac{P(A|B)P(B)}{P(A|B)P(B) + P(A|\bar{B})P(\bar{B})}.$$

Verallgemeinerung für m Ereignisse A_i ($i=1, \dots, m$):

Diese m Ereignisse **zerlegen** die Ergebnismenge, d.h. sie sind disjunkt und füllen Ω aus. Enthält Ω noch ein Ereignis B, so schneidet B mindestens ein Ereignis A_i , und B ist dann

$$P(B) = P(A_1 \cap B) + P(A_2 \cap B) + \dots + P(A_m \cap B).$$

Es gilt hier das **Bayessche Theorem**:

$$P(A_i|B) = \frac{P(A_i \cap B)}{P(B)} = \frac{P(B|A_i)P(A_i)}{\sum_{j=1}^m P(B|A_j)P(A_j)}.$$

Übung:

Was ist $P(A|B)$, falls A und B disjunkt sind?

Was ist $P(A|B)$, falls A und B stochastisch unabhängig sind?

Lösungen der Übungen

Beispiel mit den Kursverläufen

1. Darstellung der verschiedenen Wahrscheinlichkeiten

	K_1 0,55	$\overline{K_1}$ 0,45
K_2 0,55 aus Summe der Zeile	$K_1 \cap K_2$ 0,4	$\overline{K_1} \cap K_2$ 0,15
$\overline{K_2}$ 0,45 aus Summe der Zeile	$K_1 \cap \overline{K_2}$ 0,15	$\overline{K_1} \cap \overline{K_2}$ 0,3

2. Bei stochastischer Unabhängigkeit müsste die gemeinsame Wahrscheinlichkeit gleich dem Produkt der Einzelwahrscheinlichkeiten sein.

$$P(K_1 \cap K_2) = 0,4.$$

aber

$$P(K_1) \cdot P(K_2) = 0,55 \cdot 0,55 \neq 0,4.$$

Also sind die Ereignisse stochastisch abhängig.

3. Es ist

$$P(K_2|K_1) = \frac{K_1 \cap K_2}{K_1} = \frac{0,4}{0,55}$$

und

$$P(\overline{K_2}|K_1) = \frac{K_1 \cap \overline{K_2}}{K_1} = \frac{0,15}{0,55}$$

4.

$$P(K_2|\overline{K_1}) = \frac{\overline{K_1} \cap K_2}{\overline{K_1}} = \frac{0,15}{0,45}$$

Übungen zu Theoretische Erkenntnisse

Lösung: 0; $P(A)$.

1.3.

Kombinierte Zufallsvorgänge

<< hoch zum Anfang vom Inhaltsverzeichnis

<< hoch zum Kapitel 1 vom Inhaltsverzeichnis

Kombinierte Zufallsvorgänge (insbesondere **wiederholte oder mehrfache Versuche**).

Allgemeines

Beispiele für kombinierte Zufallsvorgänge:

- Eine Münze werfen, dann einmal würfeln.
- Aus einer Urne ohne Zurücklegen 3 Kugeln ziehen.

- Aus einer Lostrommel 10 Gewinner ziehen.
- Gewinnspiel: Aus drei Toren eines wählen. Falls richtiges Tor, Wahl zwischen zwei Umschlägen.
- 5x auf ein Ziel schießen.

Beispiel für die formale Definition

Es sollen nacheinander drei Zufallsexperimente durchgeführt werden. Die Wahrscheinlichkeit, dass beim ersten Versuch das Ereignis A, beim zweiten Versuch das Ereignis B und beim dritten Versuch das Ereignis C resultiert, wird bezeichnet als $P(A^{(1)} \wedge B^{(2)} \wedge C^{(3)})$. A, B und C können verschiedenen Ergebnismengen entstammen! Der hochgestellte Index kann unter Umständen weggelassen werden.

Beispiel für unabhängige Versuche

Wir betrachten den Zufallsvorgang: Wir werfen zuerst eine Münze und würfeln dann.

Die beiden Versuche haben jeweils die Ergebnismenge

$$\Omega_M = \{\text{Wappen (W); Zahl (Z)}\} \text{ bzw. } \Omega_W = \{1, 2, 3, 4, 5, 6\}$$

Es ergibt sich für diesen kombinierten Versuch die Ergebnismenge Ω^* als kartesisches Produkt von Ω_M und Ω_W :

$$\Omega^* = \{(W; 1), (W; 2), (W; 3), \dots, (W; 6), (Z; 1), (Z; 2), \dots, (Z; 6)\}.$$

Ω^* hat 12 Elemente. Jedes Element hat die selbe Wahrscheinlichkeit, gezogen zu werden.

Wir suchen nun die Wahrscheinlichkeit für das Ereignis A^* : Es wird erst Wappen geworfen und dann mindestens Fünf (F) gewürfelt:

Das Ereignis $A^* = W^{(1)} \wedge F^{(2)}$ belegt in Ω^* 2 Elemente. Wir erhalten dann für die Wahrscheinlichkeit nach dem Symmetrieprinzip

$$P(A^*) = P(W^{(1)} \wedge F^{(2)}) = \frac{2}{12} = \frac{1}{6}$$

Würfeln und Münzwurf sind jedoch stochastisch unabhängig und die Wahrscheinlichkeit muss nicht umständlich über die Ergebnismenge ermittelt werden. Also ist dann

$$P(A^*) = P(W^{(1)}) \cdot P(F^{(2)}) = \frac{1}{2} \cdot \frac{2}{6} = \frac{1}{6}$$

Übung

Sie würfeln 3 mal. Mit welcher Wahrscheinlichkeit erhalten Sie zuerst zwei mal Sechs und dann höchstens Zwei?

Lösung: $\frac{1}{108}$.

Wiederholte Versuche können aber oft stochastisch abhängig sein.

Aus einer Urne mit 2 roten und 1 schwarzen Kugeln sollen zwei Kugeln ohne Zurücklegen gezogen werden.

Das zweite Ergebnis ist vom ersten natürlich nicht mehr unabhängig, weil sich je nach erster gezogener Kugel der Inhalt der Urne ändert. Es sei: R: eine rote Kugel wird gezogen und S: eine schwarze Kugel wird gezogen.

Wir wollen zuerst die Ergebnismenge der abhängigen Versuche analysieren. Nummerieren wir die beiden roten Kugeln in R_1 und R_2 . Man kann dann bei zwei mal ziehen folgende Ergebnisse erhalten:

$$\Omega^* = \{(R_1; R_2), (R_1; S), (R_2; R_1), (R_2; S), (S; R_1), (S; R_2)\}$$

Ω^* hat insgesamt 6 Ergebnisse.

Wir definieren das Ereignis A: Zuerst wird eine rote (R), dann eine schwarze Kugel (S) gezogen, also $A = R^{(1)} \wedge S^{(2)}$.

Es gibt in Ω^* zwei Ergebnisse, die A betreffen, also ist die Wahrscheinlichkeit

$$P(A) = \frac{2}{6} = \frac{1}{3}.$$

Dieses Beispiel war einfach. Aber kann jetzt bei abhängigen Versuchen auch die Wahrscheinlichkeit für das kombinierte Ereignis unter Verzicht auf die vollständige Darstellung der Ergebnismenge bestimmt werden?

Bei stochastisch abhängigen Versuchen können die Wahrscheinlichkeiten nicht mehr ohne weiteres als Produkt der Einzelwahrscheinlichkeiten der Ereignisse bestimmt werden. Man kann aber sukzessiv den Multiplikationssatz der Ereignisse anwenden, der von den bedingten Wahrscheinlichkeiten bekannt ist: $P(A \cap B) = P(A) \cdot P(B|A)$. Die Wahrscheinlichkeit, dass beim ersten Mal A und beim zweiten Mal B resultiert, ist also

$$P(A^{(1)} \wedge B^{(2)}) = P(A^{(1)}) \cdot P(B^{(2)}|A^{(1)})$$

Es ist nach der obigen Formel

$$P(A) = P(R^{(1)} \cap S^{(2)}) = P(R^{(1)}) \cdot P(S^{(2)}|R^{(1)})$$

$$= \frac{2}{3} \cdot \frac{1}{2} = \frac{1}{3}$$

Beim ersten Versuch sind 3 Kugeln in der Urne; zwei sind rot Beim zweiten Versuch sind noch 2 Kugeln in der Urne; eine ist schwarz.

Diese Regel lässt sich auch auf mehr als zwei Ereignisse erweitern:

Beispiel

Aus einer Urne mit 10 roten (R) und 5 schwarzen (S) Kugeln sollen ohne Zurücklegen nacheinander drei rote Kugeln gezogen werden. Die Wahrscheinlichkeit dafür ist

$$P(R^{(1)} \cap R^{(2)} \cap R^{(3)}) = \frac{10}{15} \cdot \frac{9}{14} \cdot \frac{8}{13}$$

Für mehr als zwei Ereignisse kann der allgemeine **Multiplikationssatz der Wahrscheinlichkeiten** angewendet werden. Er gilt auch für Ereignisse, die nicht aus einer gemeinsamen Ergebnismenge stammen:

$$P(A^{(1)} \wedge A^{(2)} \wedge \dots \wedge A^{(m)}) = P(A^{(1)}) \cdot P(A^{(2)}|A^{(1)}) \cdot P(A^{(3)}|A^{(1)} \wedge A^{(2)}) \cdot \dots$$

$$\dots \cdot P(A^{(m)}|A^{(1)} \wedge A^{(2)} \wedge \dots \wedge A^{(m-1)}).$$

Falls die $A^{(i)}$ ($i = 1, 2, \dots, m$) stochastisch unabhängig sind, ist natürlich wieder

$$P(A^{(1)} \wedge A^{(2)} \wedge \dots \wedge A^{(m)}) = P(A^{(1)}) \cdot P(A^{(2)}) \cdot \dots \cdot P(A^{(m)}).$$

Je nachdem, wie die Problemstellung ist, gibt es für die Berechnung von Wahrscheinlichkeiten kombinierter Zufallsvorgänge also verschiedene Möglichkeiten:

1. Wir bestimmen alle Elemente von Ω^* , falls das möglich und durchführbar ist. Dann wenden wir das Symmetrieprinzip an.
2. Wir überlegen uns, beispielsweise mit Hilfe der Kombinatorik, die Zahl der Elemente in Ω^* und wenden dann das Symmetrieprinzip an.
3. Wir verwenden den allgemeinen Multiplikationssatz der Wahrscheinlichkeiten und können vielleicht sogar stochastische Unabhängigkeiten ausnützen.

Urnenmodelle

Bei wiederholten Versuchen greift man häufig auf das so genannte Urnenmodell zurück: Dieses Modell funktioniert im Prinzip folgendermaßen: Eine Urne enthält N viele Kugeln, die sich voneinander unterscheiden lassen. Es werden n viele Kugeln gezogen. Man interessiert sich für die

Zahl von Kugeln mit einem bestimmten Merkmal unter den n gezogenen.

Wir unterscheiden grundsätzlich

- das **Urnenmodell mit Zurücklegen**: Eine Kugel wird gezogen und wieder zurückgelegt
- das **Urnenmodell ohne Zurücklegen**: Eine Kugel wird gezogen und nicht wieder zurückgelegt

Viele Zufallsvorgänge, speziell die wiederholter Versuche, können auf das Urnenmodell zurückgeführt werden. Den Anfänger mag die Vorstellung, eine Kugel zu ziehen und wieder zurückzulegen, eigenartig anmuten, aber so kann man unabhängige Versuche modellieren: Betrachten wir den Zufallsvorgang, zwei mal zu würfeln, so kann man stattdessen auch aus einer Urne mit 6 verschiedenen Kugeln zwei mal jeweils eine ziehen und wieder zurücklegen.

Kombinatorik

Wir haben eine Urne mit N Kugeln gegeben. Es sollen n Kugeln gezogen werden. Wir befassen uns nun mit der Zahl der möglichen Ergebnisse bei wiederholten Versuchen. Hier müssen wir die verschiedenen Arten der Anordnung gezogener Kugeln im Urnenmodell berücksichtigen.

Zur Verdeutlichung dieser Aufgabenstellung betrachten wir eine Urne mit 3 Kugeln A, B, C. Es sollen $n = 2$ Kugeln gezogen werden. Wie viel verschiedene Paare würden wir erhalten?

Wir unterscheiden die Aufgabenstellungen

Mit Wiederholung - Mit Berücksichtigung der Reihenfolge

Die Buchstaben werden mit Zurücklegen gezogen; ein Buchstabe kann also mehrmals im Paar auftauchen. Es kommt auf die Reihenfolge der Buchstaben an. Es sind folgende verschiedene Paare möglich:

(A,A), (A,B), (A,C), (B,A), (B,B), (B,C), (C,A), (C,B), (C,C).

Es gibt insgesamt N^n viele verschiedene Ergebnisse, wie man leicht sieht.

Mit Wiederholung - Ohne Berücksichtigung der Reihenfolge

Es sind folgende verschiedene Paare möglich:

(A,A), (A,B), (A,C), (B,B), (B,C), (C,C).

Es gibt insgesamt $\binom{N+n-1}{n}$ viele verschiedene Ergebnisse.

Ohne Wiederholung - Mit Berücksichtigung der Reihenfolge

Die Buchstaben werden ohne Zurücklegen gezogen; ein Buchstabe kann nur einmal im Paar auftauchen. Es sind folgende verschiedene Paare möglich:

(A,B), (A,C), (B,A), (B,C), (C,A), (C,B).

Es gibt insgesamt $\frac{N!}{(N-n)!}$ viele verschiedene Ergebnisse.

Ohne Wiederholung - Ohne Berücksichtigung der Reihenfolge

Es sind folgende verschiedene Paare möglich:

(A,B), (A,C), (B,C).

Es gibt insgesamt $\binom{N}{n}$ viele verschiedene Ergebnisse.

Übungsbeispiel

Aus vier Personen Anna (A), Balduin (B), Cäcilie (C), Dagobert (D) werden zwei zum Geschirrspülen ausgelost, wobei eine Person abspült und eine abtrocknet.

Handelt es sich um ein Modell mit oder ohne Zurücklegen? Theoretisch wäre auch ein Modell mit Zurücklegen denkbar. Da das aber als unfair empfunden wird, gehen wir vom Modell ohne Zurücklegen (M. o. Z.) aus.

- Mit welcher Wahrscheinlichkeit erwischt es zuerst Cäcilie und dann Balduin (Ereignis E)?

Hier kommt es auf die Reihenfolge der gezogenen „Kugeln“ an.

Methode a: Direkt über die Ergebnismenge

Die Ergebnismenge ergibt $\Omega^* =$

-	(A,B)	(A,C)	(A,D)
(B,A)	-	(B,C)	(B,D)
(C,A)	(C,B)	-	(C,D)
(D,A)	(D,B)	(D,C)	-

Jedes Paar hat die gleiche Wahrscheinlichkeit, gewählt zu werden. Es gibt insgesamt $|\Omega^*| = 12$ verschiedene Paare.

$$P(E) = P((C, B)) = \frac{1}{12}$$

Methode b: Über die Zahl der Ergebnisse Es handelt sich um ein Modell ohne Zurücklegen mit Beachtung der Reihenfolge. Es gibt

$$\frac{N!}{(N-n)!} = \frac{4!}{(4-2)!} = \frac{1 \cdot 2 \cdot 3 \cdot 4}{2} = 12$$

verschiedene Paare. Es gibt nur ein Ergebnis für das Ereignis E. Es ist also

$$P(E) = \frac{|E|}{|\Omega^*|} = \frac{1}{12}$$

Methode c: Über den Multiplikationssatz der Wahrscheinlichkeiten

$$P(C^{(1)} \cap B^{(2)}) = \frac{1}{4} \cdot \frac{1}{3} = \frac{1}{12}$$

- Mit welcher Wahrscheinlichkeit müssen die zwei Männer abwaschen (Ereignis F)?

Methode a:

Es ist $F = \{(B,D), (D,B)\}$. Dieses Ereignis belegt in Ω^* zwei Elemente. Also ist

$$P(F) = \frac{2}{12} = \frac{1}{6}$$

Methode b:

M.o.Z, ohne Beachtung der Reihenfolge. Es gibt

$$\binom{N}{n} = \frac{N!}{n!(N-n)!} = \frac{1 \cdot 2 \cdot 3 \cdot 4}{(1 \cdot 2)(1 \cdot 2)} = 6$$

verschiedene Paare . Es ist also $P(F) = \frac{1}{6}$

Methode c:

$$P(F) = \frac{2}{4} \cdot \frac{1}{3} = \frac{1}{6}$$

2.	Zufallsvariablen	<< hoch zum Anfang vom Inhaltsverzeichnis << hoch zum Kapitel 2 vom Inhaltsverzeichnis
----	------------------	---

Beispiel zum Begriff der Zufallsvariablen

Die fränkische Druckerei Printzig nennt 10 multifunktionelle Hochleistungsdrucker ihr eigen. Drei Drucker sind von der Firma Alpha, zwei sind von Beta, vier von Gamma und einer stammt von der Firma Delta. Da die Drucker auch von Kunden bedient werden, fallen sie aufgrund unsachgemäßer Handhabung häufig aus. Man hat festgestellt, dass alle Drucker in gleichem Maße anfällig sind. Wegen der Gewährleistung wird bei jedem Ausfall ein Wartungstechniker der betreffenden Firma geholt. Die Kosten für die Wiederherstellung eines Druckers hängen vom Hersteller ab, wobei die Drucker der Firma Gamma in der Reparatur am billigsten sind.

Am liebsten ist es natürlich Herrn Printzig, wenn ein Drucker mit den geringsten Reparaturkosten ausfällt.

Überlegen wir:

Welche Ergebnismenge gehört zu dem Zufallsvorgang: Ein Drucker fällt zufällig aus?

Mit welcher Wahrscheinlichkeit entstehen Herrn Printzig die geringsten Kosten?

Wir erhalten die Ergebnismenge

$$\Omega = \{A_1, A_2, A_3, B_1, B_2, G_1, G_2, G_3, G_4, D_1\},$$

wobei z.B. B_2 Drucker Nr. 2 der Firma Beta bedeutet. G sei das Ereignis, die geringsten Reparaturkosten zu haben. Jeder Drucker hat die gleiche Wahrscheinlichkeit, auszufallen. Dann ist nach dem Symmetrieprinzip

$$P(G) = \frac{\text{Zahl der G-Drucker}}{\text{Zahl aller Drucker}} = \frac{|G|}{|\Omega|} = \frac{4}{10} = 0,4$$

Die Kosten für die Reparatur eines Druckers betragen je nach Hersteller wie folgt:

Hersteller	Alpha	Beta	Gamma	Delta
Kosten (Euro)	50	60	30	100

Überlegen wir: Wieviel muss Herr Printzig pro Ausfall im Durchschnitt bezahlen?

Ordnen wir nun der Ergebnismenge die entsprechenden Kosten zu:

A ₁	A ₂	A ₃	B ₁	B ₂	G ₁	G ₂	G ₃	G ₄	D ₁
↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
50	50	50	60	60	30	30	30	30	100

Ω hat 10 Ergebnisse und jedes Elementarereignis hat die Wahrscheinlichkeit 1/10. Jeder Drucker fällt dann auch mit der Wahrscheinlichkeit 1/10 aus. Die durchschnittlichen Reparaturkosten sind also

$$\begin{aligned}
& 50 \cdot \frac{1}{10} + 50 \cdot \frac{1}{10} + 50 \cdot \frac{1}{10} + 60 \cdot \frac{1}{10} + 60 \cdot \frac{1}{10} + \dots + 100 \cdot \frac{1}{10} \\
&= 50 \cdot \frac{3}{10} + 60 \cdot \frac{2}{10} + 30 \cdot \frac{4}{10} + 100 \cdot \frac{1}{10} \\
&= \frac{150}{10} + \frac{120}{10} + \frac{120}{10} + \frac{100}{10} = \frac{490}{10} = 49 \text{ Euro}
\end{aligned}$$

Wir haben soeben eine Zufallsvariable konstruiert und zwar, indem wir allen Ergebnissen von Ω eine Zahl zugeordnet haben.

Den Durchschnitt konnten wir erst berechnen, nachdem wir die Drucker mit einer Zahl versehen hatten. Man kann je nach Interesse den Elementarereignissen beliebige Zahlen zuordnen. So könnten für die laufende Wartung wieder ganz andere Kosten gelten. Nur die Ergebnismenge ist festgelegt. Man könnte nun die Wahrscheinlichkeit berechnen, dass bei einem Ausfall 60 Euro fällig werden: Es gibt 10 Elementarereignisse und zwei davon entsprechen 60 Euro. Also beträgt diese Wahrscheinlichkeit $2/10$.

Wir bezeichnen eine Zufallsvariable mit einem großen Buchstaben. Die Werte, die eine Zufallsvariable annehmen kann, nennt man Ausprägung. Eine bestimmte Ausprägung kennzeichnen wir mit einem Kleinbuchstaben. Nennen wir unsere Zufallsvariable „Reparaturkosten“ X . Wir fassen jetzt die verschiedenen Wahrscheinlichkeiten der Zufallsvariablen X in einer Wahrscheinlichkeitstabelle zusammen. Herr Printzig hat 4 mal die „Chance“, 30 Euro zu bezahlen, also ist die Wahrscheinlichkeit, dass $X = 30$ ist, gleich $4/10$, usw.

Wahrscheinlichkeitstabelle:

	x_1	x_2	x_3	x_4
Ausprägung x_i	30	50	60	100
Wahrscheinlichkeit $f(x_i)$	0,4	0,3	0,2	0,1

$f(x)$ bezeichnet die zur bestimmten Ausprägung x gehörende Wahrscheinlichkeit. Es ist beispielsweise

$$P(X = 60) = f(x_3) = f(60) = 0,2,$$

aber

$$P(X = 70) = f(70) = 0,$$

denn für $X = 70$ existiert kein Ergebnis.

Die Summe aller Wahrscheinlichkeiten ist

$$\sum_{i=1}^m f(x_i) = 1$$

Man kann diese Wahrscheinlichkeiten auch grafisch als Stabdiagramm darstellen.

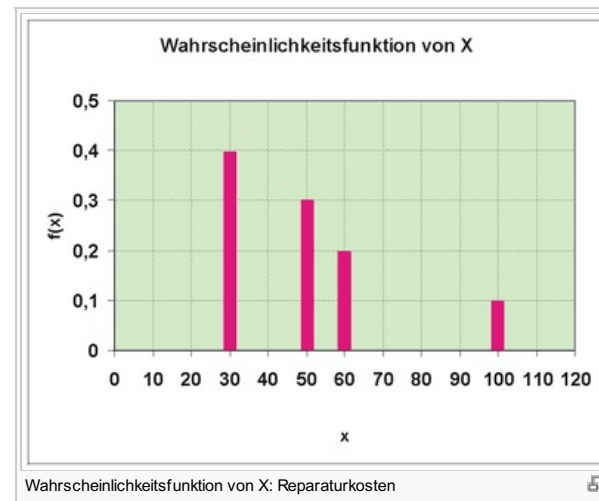
Man sieht, dass an den x -Stellen 30, 50, 60 und 100 die Wahrscheinlichkeitsfunktion die Werte 0,4, 0,3, 0,2 und 0,1 annimmt, aber an allen sonstigen Werten von x Null ist.

Wie groß ist nun aber die Wahrscheinlichkeit, dass Herr Printzig höchstens 50 Euro bezahlen muss?

$$P(X \leq 50) = P(X = 30) + P(X = 50) = 0,4 + 0,3 = 0,7.$$

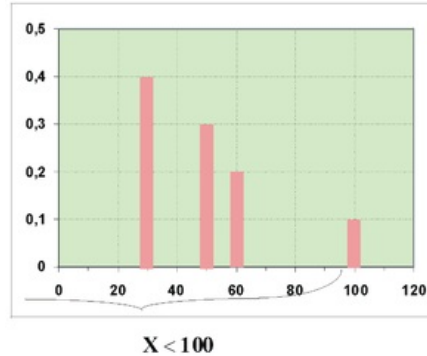
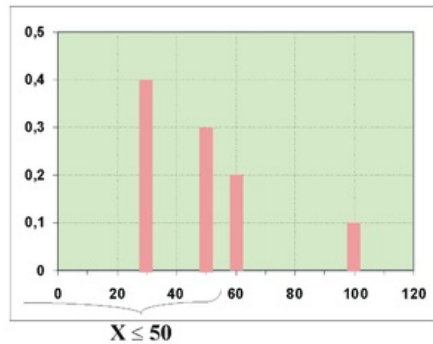
Das kann man auch aus der Graphik erkennen: Es ist die Summe der „Stäbchen“ für $x \leq 50$.

Mit welcher Wahrscheinlichkeit muss Herr Printzig weniger als 100 Euro zahlen? Gefragt ist hier nach $P(X < 100)$. Ein Blick auf die Grafik verrät



uns, dass gilt

$$P(X < 100) = P(X \leq 60) = P(X = 30) + P(X = 50) + P(X = 60) = 0,4 + 0,3 + 0,2 = 0,9.$$



Wieviel ist nun $P(30 < X \leq 60)$?

Man kann hier wieder die „Stäbchenmethode“ anwenden:

$$P(30 < X \leq 60) = 0,3 + 0,2 = 0,5.$$

Es gibt aber auch eine Rechenregel, die man mit Hilfe der Grafik leicht erkennt:

$$P(a < X \leq b) = P(X \leq b) - P(X \leq a),$$

also

$$P(30 < X \leq 60) = P(X \leq 60) - P(X \leq 30) = 0,9 - 0,4 = 0,5.$$

Die Wahrscheinlichkeiten $P(X \leq a)$ einer bestimmten Ausprägung a von X bilden die Verteilungsfunktion von X , die die Wahrscheinlichkeitsverteilung von X in eindeutiger Weise beschreibt. Das ist eine Festlegung, die die Statistiker als sinnvoll erachten. Die Verteilungsfunktionen werden grossbuchstabig als $F(a)$ bezeichnet. Meist wird statt a das Symbol x verwendet. Wir wollen die Verteilungsfunktion konstruieren, indem wir die obige Graphik zu Hilfe nehmen und für einzelne Stützwerte x die Verteilungsfunktion berechnen.

Wie groß ist z.B. $P(X \leq 10)$? Es ist $P(X \leq 10) = F(10) = 0$.

Ebenso sind $P(X \leq 15) = 0$ und $P(X \leq 20) = 0$.

Es ist also $F(a) = 0$ für alle Werte von a mit $-\infty < a < 30$.

Als nächstes untersuchen wir $P(X \leq 30)$:

$P(X \leq 30) = F(30) = 0,4$. Ebenso sind $P(X \leq 30,1) = 0,4$ und $P(X \leq 49,99999) = 0,4$.

Die Verteilungsfunktion hat also den Wert $F(a) = 0,4$ für $30 \leq a < 50$.

Es gilt weiter: $P(X \leq 50)$, $P(X \leq 59)$, ... $P(X < 60)$ sind, siehe Graphik: $0,4 + 0,3 = 0,7$.

...

Schließlich ist die Wahrscheinlichkeit $P(X \leq 100)$ oder auch $P(X \leq 110)$, $P(X \leq 1000)$ usw... gleich 1.

Wir können die Wahrscheinlichkeiten zusammenfassen in der Verteilungsfunktion

$$F(x) = P(X \leq x)$$

$$P(X \leq a) = F(a) = \begin{cases} 0 & \text{für } a < 30 \\ 0,4 & \text{für } 30 \leq a < 50 \\ 0,7 & \text{für } 50 \leq a < 60 \\ 0,9 & \text{für } 60 \leq a < 100 \\ 1 & \text{für } a \geq 100 \end{cases}$$

Man sieht, dass diese Verteilungsfunktion grafisch eine Treppenfunktion darstellt. Die Punkte links an den Stufen zeigen an, dass der Funktionswert dieser Stufe genau zum Punkt a gehört.

Man kann hier auch die Wahrscheinlichkeiten der Grafik entnehmen, z.B. ist $P(X \leq 70) = 0,9$.

Besonders interessiert man sich bei einer Zufallsvariable für zwei Kennwerte, Parameter genannt, die die Zufallsvariable genauer beschreiben.

Einer ist der durchschnittliche Wert, den die Zufallsvariable „auf lange Sicht“ annimmt, wenn der Zufallsvorgang „sehr oft“ durchgeführt wird. Dieser Parameter wird Erwartungswert EX genannt, also der Wert, den man langfristig erwarten kann. Wir hatten ihn schon oben ermittelt als

$$EX = 50 \cdot \frac{3}{10} + 60 \cdot \frac{2}{10} + 30 \cdot \frac{4}{10} + 100 \cdot \frac{1}{10} = 49$$

die durchschnittlichen Reparaturkosten.

Ein weiterer Parameter ist die Streuung der X , ein Maß, wie stark die einzelnen Werte von X von EX abweichen, also 30-49, 50-49, 60-49, 100-49. Da z.B. 100 viel seltener auftritt als 30, gewichtet man auch diese Abweichungen mit ihrer Wahrscheinlichkeit. Eine Quadrierung sorgt dann einerseits dafür, dass sich positive und negative Abweichungen nicht aufheben, andererseits für eine überproportionale Berücksichtigung von besonders starken Abweichungen. Man erhält im Ergebnis als durchschnittliche quadratische Abweichung der X -Werte von EX die Varianz

$$\begin{aligned} \text{Var } X &= (30 - 49)^2 \cdot 0,4 + (50 - 49)^2 \cdot 0,3 \\ &\quad + (60 - 49)^2 \cdot 0,2 + (100 - 49)^2 \cdot 0,1 \\ &= 361 \cdot 0,4 + 1 \cdot 0,3 + 121 \cdot 0,2 + 2601 \cdot 0,1 = 429 \end{aligned}$$

wobei zu beachten ist, dass sich hier als Einheit Euro² ergibt.

Die Wurzel der Varianz ist die Standardabweichung; man könnte sie salopp als mittlere Abweichung der Ausprägungen vom Durchschnitt bezeichnen. Sie beträgt in unserem Beispiel etwa 20,71.

Allgemeine Darstellung einer Zufallsvariablen

Gegeben ist ein Zufallsvorgang mit der Ergebnismenge Ω . Jedem Element aus Ω wird eine reelle Zahl x zugeordnet:

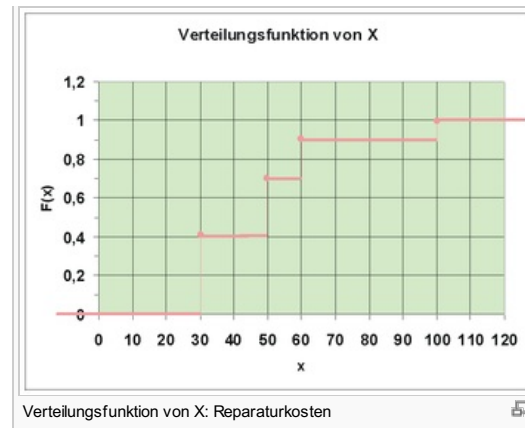
$$\Omega \rightarrow \mathbb{R}, \quad \omega \mapsto X(\omega).$$

Die Elemente von X sind Realisationen, Ausprägungen, Werte. Die Verteilung der Zufallsvariablen kann festgelegt werden mit ihrer Verteilungsfunktion F , definiert als

$$F(x) = P(X \leq x).$$

Es gilt für die Verteilung jeder Zufallsvariablen:

- $F(x)$ ist für alle $x \in \mathbb{R}$ definiert.
- $0 \leq F(x) \leq 1$.
- $F(x)$ ist monoton steigend, also $x_1 < x_2 \rightarrow F(x_1) \leq F(x_2)$
- $F(x)$ ist rechtsseitig stetig.
- $P(a < X \leq b) = P(X \leq b) - P(X \leq a)$.



2.1.

Diskrete Zufallsvariablen

[<< hoch zum Anfang vom Inhaltsverzeichnis](#)

[<< hoch zum Kapitel 2 vom Inhaltsverzeichnis](#)

Eine Zufallsvariable ist diskret, wenn sie in jedem beschränkten Intervall nur endlich viele Ausprägungen annehmen kann. Die diskrete Zufallsvariable kann endlich oder abzählbar unendlich viele Werte x_i ($i = 1, 2, \dots, m$ bzw. $i = 1, 2, \dots$) annehmen.

Beispiele

- Zahl der Schadensleistungen, die in einem Jahr bei einer Versicherung auftreten
- Kinderzahl von Konsumenten
- Zahl der defekten Kondensatoren in einem Fertigungslos

Ihre Wahrscheinlichkeitsfunktion ist

$$P(X = x) = f(x) = \begin{cases} f(x_i) & \text{für } x = x_i \\ 0 & \text{sonst} \end{cases}$$

Es gilt

$$\sum_i f(x_i) = 1.$$

Die Verteilungsfunktion $P(X \leq a) = F(a)$ ist die Summe aller Wahrscheinlichkeiten $f(x_i)$ für $x_i \leq a$.

Der Erwartungswert einer Zufallsvariablen ist der Durchschnitt des Auftretens ihrer Realisationen. Bei einer diskreten Zufallsvariablen beträgt er

$$EX = \sum_i x_i f(x_i),$$

falls EX existiert, d.h. nicht unendlich wird.

Die Varianz einer diskreten Zufallsvariablen berechnet sich als

$$\text{Var } X = \sum_i (x_i - EX)^2 f(x_i).$$

Nach dem sog. Verschiebungssatz ist auch

$$\text{Var } X = \left(\sum_i x_i^2 f(x_i) \right) - (EX)^2,$$

im Beispiel:

$$\begin{aligned} \text{Var } X &= 30^2 \cdot 0,4 + 50^2 \cdot 0,3 + 60^2 \cdot 0,2 + 100^2 \cdot 0,1 - 49^2 \\ &= 360 + 750 + 720 + 1000 - 2401 = 429. \end{aligned}$$

2.2.

Stetige Zufallsvariablen

[<< hoch zum Anfang vom Inhaltsverzeichnis](#)

[<< hoch zum Kapitel 2 vom Inhaltsverzeichnis](#)

Beispiel eines Zeitungskiosks

Dichtefunktion

Die Zufallsvariable X : „An einem Tag verkaufte Menge an Tageszeitungen (in 100) eines Zeitungskiosks“ lässt sich beschreiben mit der (in diesem Fall frei erfundenen) **Dichtefunktion**

Dichtefunktion

$$f(x) = \begin{cases} \frac{1}{4}x - \frac{3}{2} & \text{für } 6 \leq x \leq 8 \\ \frac{5}{2} - \frac{1}{4}x & \text{für } 8 < x \leq 10 \\ 0 & \text{sonst} \end{cases}$$

Diese Zufallsvariable X ist nun **stetig**, d.h. sie hat in jedem Intervall $a \leq X \leq b$ unendlich viele Ausprägungen.

Eine Analyse der Grafik zeigt, dass diese Dichtefunktion symmetrisch bezüglich 8 ist, was die Berechnung von Wahrscheinlichkeiten sehr erleichtert.

Wir wollen nun die Wahrscheinlichkeit bestimmen, dass an einem Tag höchstens 700 Zeitungen verkauft werden, also $P(X \leq 7)$. Wenn wir analog zu der diskreten Zufallsvariablen vorgehen, wo wir „die Summe der Stäbchen“ ermittelten, müsste die Wahrscheinlichkeit $P(X \leq a)$ hier „unendlich viele Stäbchen“, also eine **Fläche** ergeben.

Wir berechnen die Dreiecksfläche mit Hilfe der Geometrie:

$$\begin{aligned} P(X \leq 7) &= \text{Breite des Dreiecks} \cdot \text{Höhe des Dreiecks} \cdot \frac{1}{2} \\ &= 1 \cdot \frac{1}{4} \cdot \frac{1}{2} = \frac{1}{8} . \end{aligned}$$

Es ist übrigens auch

$$P(X < 7) = \frac{1}{8} ,$$

denn bei einer stetigen Zufallsvariablen ist $P(X = x) = 0$, da es als unmöglich angesehen wird, genau einen bestimmten Wert x zu „treffen“. Man betrachtet also bei einer stetigen Zufallsvariablen nur Wahrscheinlichkeiten der Art $P(X \leq x)$ o.ä.

Es ist $P(X \leq 8) = 0,5$, wie man der Grafik sofort entnimmt.

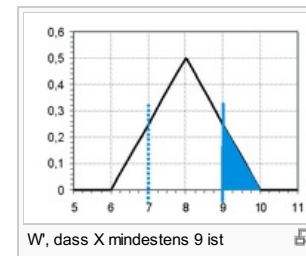
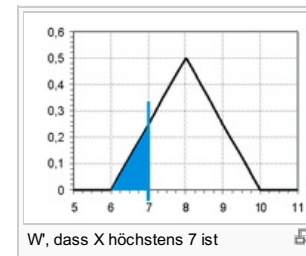
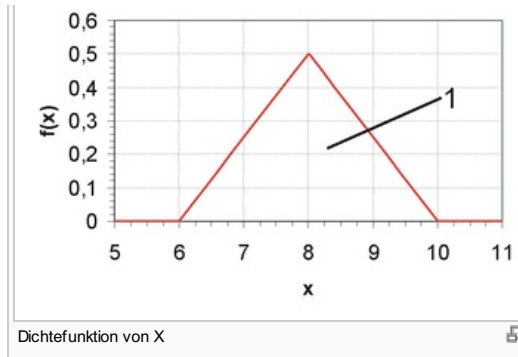
$P(X \geq 9) = \frac{1}{8}$, denn wie man sieht, ist die Fläche von $P(X \geq 9)$ genau gleich der Fläche $P(X \leq 7)$.

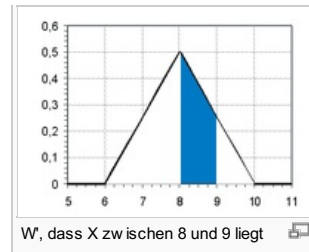
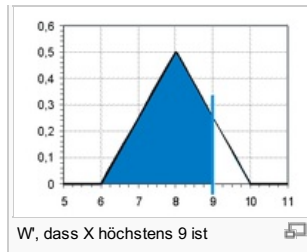
Außerdem ist $P(X \leq 9) = 1 - P(X \geq 9) = \frac{7}{8}$.

Bestimmen wir die Wahrscheinlichkeit eines Intervalls. Es ergibt

$$P(8 < X \leq 9) = P(X \leq 9) - P(X \leq 8) = 0,875 - 0,5 = 0,375,$$

wenn man die Rechenregel für $P(a < X \leq b)$ anwendet.





Verteilungsfunktion

Man kann Wahrscheinlichkeiten von X auch als **Verteilungsfunktion** darstellen. Sucht man die Wahrscheinlichkeit $P(X \leq a)$, muss also das Integral von $-\infty$ bis a berechnet werden:

$$P(X \leq a) = F(a) = \int_{-\infty}^a f(x) dx$$

Bei unserem Beispiel sind wir mit verschiedenen Bereichen konfrontiert:

1. $a < 6$

$$P(X \leq a) = F(a) = \int_{-\infty}^a 0 dx = 0$$

2. $6 \leq a \leq 8$

$$\begin{aligned} F(a) &= \int_{-\infty}^6 0 dx + \int_6^a \left(\frac{1}{4}x - \frac{3}{2} \right) dx \\ &= 0 + \left[\frac{x^2}{8} - \frac{3}{2}x \right]_6^a \\ &= \frac{a^2}{8} - \frac{3}{2}a - \left(\frac{6^2}{8} - \frac{3}{2} \cdot 6 \right) = \frac{a^2}{8} - \frac{3}{2}a + \frac{9}{2} \end{aligned}$$

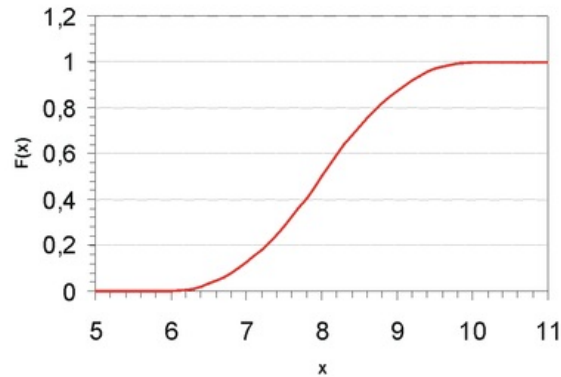
3. $8 < a \leq 10$

$$\begin{aligned} F(a) &= \int_{-\infty}^6 0 dx + \int_6^8 \left(\frac{1}{4} \cdot x - \frac{3}{2} \right) dx + \int_8^a \left(\frac{5}{2} - \frac{1}{4}x \right) dx \\ &= 0 + \left[\frac{x^2}{8} - \frac{3}{2} \cdot x \right]_6^8 + \left[\frac{5}{2} \cdot x - \frac{x^2}{8} \right]_8^a \\ &= \left(\frac{64}{8} - \frac{3}{2} \cdot 8 \right) - \left(\frac{36}{8} - \frac{3}{2} \cdot 6 \right) + \left(\frac{5}{2} \cdot a - \frac{a^2}{8} \right) - \left(\frac{5}{2} \cdot 8 - \frac{64}{8} \right) \\ &= -\frac{a^2}{8} + \frac{5}{2} \cdot a - \frac{23}{2} \end{aligned}$$

4. $a > 10$

$$F(a) = 1$$

Verteilungsfunktion



$$P(X \leq a) = F(a) = \begin{cases} 0 & \text{für } a < 6 \\ \frac{a^2}{8} - \frac{3}{2}a + \frac{9}{2} & \text{für } 6 \leq a \leq 8 \\ -\frac{a^2}{8} + \frac{5}{2}a - \frac{23}{2} & \text{für } 8 < a \leq 10 \\ 1 & \text{sonst} \end{cases}$$

Wir erhalten beispielsweise durch Einsetzen in $F(x)$

$$P(X \leq 7) = F(7) = \frac{7^2}{8} - \frac{3}{2} \cdot 7 + \frac{9}{2} = \frac{1}{8},$$

$$P(X \leq 9) = F(9) = -\frac{9^2}{8} + \frac{5}{2} \cdot 9 - \frac{23}{2} = \frac{7}{8}.$$

Quantil

Das Quantil $x(p)$ gibt die Ausprägung x an, die zu einem bestimmten Verteilungswert $p = F(x)$ gehört.

Beispiele

$x(0,875) = 9$, d.h. zur Wahrscheinlichkeit 0,875 gehört der x -Wert 9.

Ebenso ist $x(0,5) = 8$. D.h. 8 ist der Median, also wurden an 50% aller Tage höchstens 800 Zeitungen verkauft.

Übung

Bestimmen Sie $P(6,25 < X < 8,75)$. Mit welcher Wahrscheinlichkeit wurden an den 50% besten Tagen mindestens 900 Zeitungen verkauft? Gesucht ist hier $P(X > 9 | X > 8)$.

Was Sie speziell über stetige Zufallsvariablen wissen sollten

Eine stetige Zufallsvariable kann in jedem beschränkten Intervall unendlich viele Ausprägungen annehmen. Ihre Verteilung lässt sich durch eine Dichtefunktion $f(x)$ beschreiben. $f(x)$ ist keine Wahrscheinlichkeit, sondern eine Dichte.

- Die Verteilungsfunktion ist

$$P(X \leq a) = F(a) = \int_{-\infty}^a f(x) dx$$

- Es gilt: $P(X = a) = 0$.

- Wegen $P(X = a) = 0$ ist $P(X \leq a) = P(X < a)$ und $P(X > a) = P(X \geq a)$
- Die Dichtefunktion $f(x)$ ist die erste Ableitung der Verteilungsfunktion, falls diese an der Stelle x differenzierbar ist.
- Die Dichtefunktion $f(a)$ kann auch größer als 1 werden.
- Ausgehend von $P(X \leq x) = p$ ist das p -Quantil $x(p)$ der Wert x , der zu einer gegebenen Wahrscheinlichkeit p gehört. Speziell $x(0,5)$ ist der Median.
- Der Erwartungswert einer stetigen Zufallsvariablen ist analog zu oben

$$EX = \int_{-\infty}^{\infty} x \cdot f(x) dx, \text{ falls } EX \text{ existiert, d.h. nicht unendlich wird.}$$

- Ihre Varianz ist

$$\text{Var } X = \int_{-\infty}^{\infty} (x - EX)^2 \cdot f(x) dx$$

wobei auch hier der Verschiebungssatz angewendet werden kann:

$$\text{Var } X = \left(\int_{-\infty}^{\infty} x^2 f(x) dx \right) - (EX)^2$$

Bei symmetrisch verteilten Zufallsvariablen ist im Allgemeinen der Erwartungswert der Zufallsvariablen gleich dem Median.

In unserem Beispiel ist also $EX = 8$, denn die Verteilung ist symmetrisch. Das bedeutet, dass im Durchschnitt pro Tag 800 Zeitungen umgesetzt werden.

Wendet man die gegebene Formel für EX auf unser Beispiel an, so erhält man:

$$\begin{aligned} EX &= \int_{-\infty}^6 x \cdot 0 dx + \int_6^8 x \cdot \left(\frac{x}{4} - \frac{3}{2}\right) dx + \int_8^{10} x \cdot \left(\frac{5}{2} - \frac{x}{4}\right) dx + \int_{10}^{\infty} x \cdot 0 dx \\ &= \left[\frac{x^3}{12} - \frac{3x^2}{4} \right]_6^8 + \left[\frac{5x^2}{4} - \frac{x^3}{12} \right]_8^{10} = 8 \end{aligned}$$

Entsprechend gilt für die Varianz:

$$\begin{aligned} \text{Var } X &= \int_{-\infty}^6 x^2 \cdot 0 dx + \int_6^8 x^2 \cdot \left(\frac{x}{4} - \frac{3}{2}\right) dx + \int_8^{10} x^2 \cdot \left(\frac{5}{2} - \frac{x}{4}\right) dx \\ &\quad + \int_{10}^{\infty} x^2 \cdot 0 dx - 64 \\ &= \left[\frac{x^4}{16} - \frac{3x^3}{6} \right]_6^8 + \left[\frac{5x^3}{6} - \frac{x^4}{16} \right]_8^{10} - 64 = \frac{2}{3} \approx 0,7 \end{aligned}$$

Beispiel: Eingehende Anrufe bei Fernsehabstimmungen

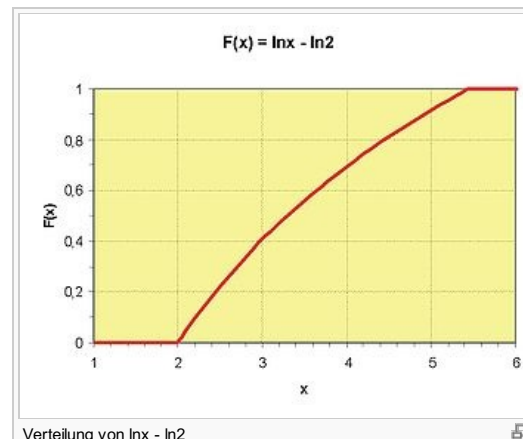
Während einer Fernsehsendung wurden die Zuschauer aufgefordert, telefonisch abzustimmen. Die Leitungen wurden um 14 Uhr freigeschaltet. Dann konnten die Zuschauer bis ca. 17.30 Uhr anrufen. Für die eintreffenden Anrufe ergab sich näherungsweise die Verteilungsfunktion der stetigen Zufallsvariablen X : Zeitpunkt (Uhrzeit), an dem ein Anruf eintrifft, wie folgt:

$$F(x) = \begin{cases} 0 & \text{für } x < 2 \\ \ln x - \ln 2 & \text{für } 2 \leq x \leq 2e \\ 1 & \text{für } x > 2e \end{cases}$$

Sei jetzt $\omega \in \Omega$ ein beliebiger Anruf.

Wir wollen nun bestimmen

1. die Dichtefunktion $f(x)$
2. die Wahrscheinlichkeit dass bis höchstens 15 Uhr der Anruf ω eingegangen ist.



3. die Wahrscheinlichkeit, dass zwischen 15 und 16 Uhr der Anruf ω eingegangen ist.
4. die Uhrzeit, zu der 90% aller Anrufe eingetroffen sind
5. den Median
6. den Erwartungswert
7. die Varianz

Die Grafik der Verteilung $F(X)$ zeigt den typischen Verlauf einer logarithmischen Funktion.

1. Dichtefunktion $f(x)$

Die Dichtefunktion ist immer die erste Ableitung der Verteilungsfunktion: $f(x) = F'(x)$.

Unsere Verteilungsfunktion ist abschnittsweise definiert. Wir müssen bereichsweise ableiten (dass die Funktion an den Knickstellen möglicherweise nicht differenzierbar ist, tut im Allgemeinen nicht weh, Hauptsache, die Fläche ergibt 1).

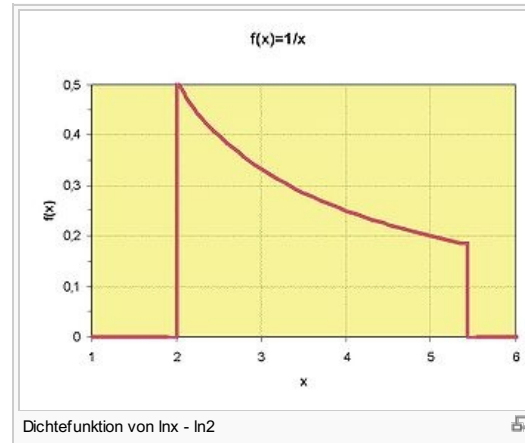
$$\text{Bereich } x < 2: F(x) = 0 \rightarrow f(x) = 0$$

$$\text{Bereich } 2 \leq x \leq 2e: F(x) = \ln x - \ln 2 \rightarrow f(x) = \frac{1}{x}$$

$$\text{Bereich } x > 2e: F(x) = 1 \rightarrow f(x) = 0$$

Wir wollen jetzt $f(x)$ noch ordentlich angeben:

$$f(x) = \begin{cases} \frac{1}{x} & \text{für } 2 \leq x \leq 2e \\ 0 & \text{sonst} \end{cases}$$



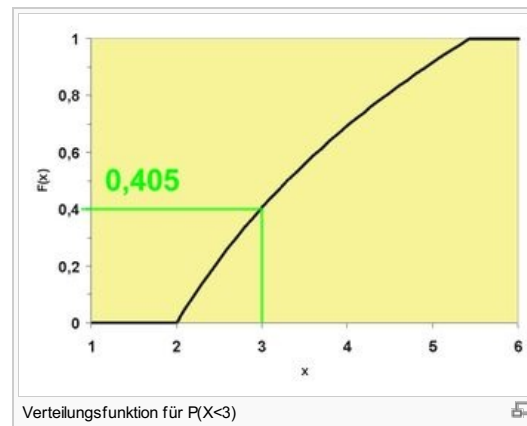
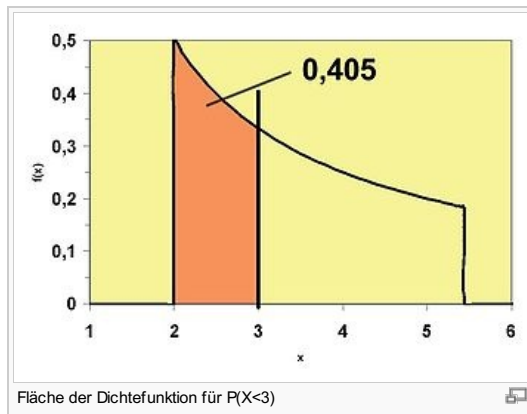
Betrachten wir mal die Dichtefunktion: Man sieht hier deutlich, dass die meisten Anrufe in den ersten 1,5 Stunden nach Freischalten eingelaufen sind. Danach flaut die Zahl der Anrufe allmählich ab.

2. Wahrscheinlichkeit, dass bis höchstens 15 Uhr der Anruf ω eingegangen ist

Gesucht ist $P(X \leq 3)$. In der Dichtefunktion ist das die Fläche von 2 bis 3. Diese Fläche ist das Selbe wie der Funktionswert $F(3)$. Wir erhalten

$$P(X \leq 3) = \ln 3 - \ln 2 = 1,0986 - 0,6931 = 0,4055$$

Man kann also sagen, dass in einer Stunde ab Freischalten der Leitungen 40% der Anrufe eingegangen sind.



3. Wahrscheinlichkeit, dass zwischen 15 und 16 Uhr der Anruf w eingegangen ist

Gesucht ist hier $P(3 \leq X \leq 4)$. Wir wissen schon, dass speziell für stetige Zufallsvariablen (bei diskreten muss man noch zwischen $<$ und \leq unterscheiden) gilt: $P(3 \leq X \leq 4) = P(X \leq 4) - P(X \leq 3)$. Wir erhalten dann

$$\begin{aligned} P(3 \leq X \leq 4) &= F(4) - F(3) \\ &= \ln 4 - \ln 2 - (\ln 3 - \ln 2) \\ &= \ln 4 - \ln 3 \\ &= 1,3863 - 1,0986 = 0,2877 \end{aligned}$$

4. Uhrzeit, zu der 90% aller Anrufe eingetroffen sind

Hier ist die Wahrscheinlichkeit 0,9 gegeben und wir suchen den X-Wert, der zu dieser Wahrscheinlichkeit passt, also $P(X \leq ?) = 0,9$. Gesucht ist also das 90%-Quantil. Wir bilden den Ansatz

$F(?) = 0,9$ oder etwas professioneller: $F(x(0,9)) = 0,9$, also

$$\ln x - \ln 2 = 0,9 \rightarrow \ln x = \ln 2 + 0,9 \rightarrow x = \exp(\ln 2 + 0,9) \approx 4,91.$$

d.h. etwa um 16.55 waren 90% der Anrufe eingegangen.

5. Median

Der Median ist das 50%-Quantil. Es ergibt sich also analog zu oben:

$$\ln x - \ln 2 = 0,5 \rightarrow x \approx 3,30$$

6. Erwartungswert

Der Erwartungswert der Zufallsvariablen X wird bei einer stetigen Zufallsvariablen integriert:

$$EX = \int_{-\infty}^{\infty} x \cdot f(x) dx$$

Wir müssen hier wieder bereichsweise vorgehen und bestimmen zunächst mal die Teilintegrale:

$$\text{Bereich } x < 2: \int_{-\infty}^2 x \cdot 0 dx = 0$$

$$\text{Bereich } 2 \leq x \leq 2e: \int_2^{2e} x \cdot \frac{1}{x} dx = \int_2^{2e} 1 dx = [x]_2^{2e} = 2e - 2 = 3,44.$$

$$\text{Bereich } x > 2e: \int_{2e}^{\infty} x \cdot 0 dx = 0$$

Wir müssen nun die Teilintegrale addieren und erhalten

$$EX = 0 + 3,44 + 0 = 3,44.$$

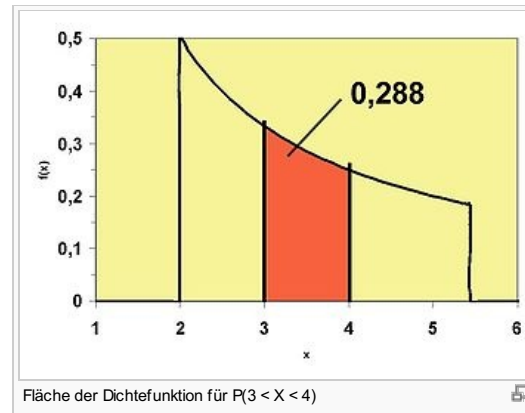
Es kam also ein Anruf im Durchschnitt um 15.30 an.

7. Varianz

Die Varianz berechnet sich nach der Formel

$$\text{Var } X = \left(\int_{-\infty}^{\infty} x^2 \cdot f(x) dx \right) - (EX)^2.$$

Analog zu oben erhalten wir



$$= \left(\int_2^{2e} x^2 \cdot \frac{1}{x} dx \right) - (EX)^2 = \left(\int_2^{2e} x dx \right) - (2e - 2)^2$$

$$= \left[\frac{x^2}{2} \right]_2^{2e} - 3,44^2 = \frac{(2e)^2}{2} - \frac{2^2}{2} - 3,44^2 = 0,9681$$

2.3.

Ungleichung von Bienaymé-Tschebyschew

<< hoch zum Anfang vom Inhaltsverzeichnis

<< hoch zum Kapitel 2 vom Inhaltsverzeichnis

Mit der [Ungleichung von Tschebyschew](#) oder Bienaymé-Tschebyschew kann man Wahrscheinlichkeiten einer Zufallsvariablen mit unbekannter Verteilung abschätzen. Benötigt werden als Information der Erwartungswert und die Varianz der Zufallsvariablen, die im Allgemeinen geschätzt werden müssen.

Die Ungleichung lautet folgendermaßen:

$$P(|X - EX| \geq \epsilon) \leq \frac{\text{Var}X}{\epsilon^2}.$$

Besser kann man sich die Beziehung vorstellen, wenn man die Betragsungleichung ausschreibt :

$$P(X \leq EX - \epsilon \vee X \geq EX + \epsilon) \leq \frac{\text{Var}X}{\epsilon^2}$$

Diese Abschätzung ist naturgemäß sehr grob und kann manchmal nichtssagende Ergebnisse liefern.

Beispiel

Es ist bekannt, dass ein Kaffeeautomat im Durchschnitt 250 ml Kaffee ausschenkt mit einer Varianz von 100 ml². Eine Tasse gilt als korrekt befüllt, wenn ihr Inhalt nicht mehr als 30 ml vom Durchschnitt abweicht. Der Anteil der inkorrekt befüllten Tassen beträgt höchstens

$$P(|X - 250| \geq 30) \leq \frac{100}{30^2} = \frac{1}{9}.$$

bzw.

$$P(X \leq EX - 30 \vee X \geq EX + 30) \leq \frac{100}{30^2} = \frac{1}{9}.$$

Umgekehrt gilt dann auch

$$P(EX - \epsilon < X < EX + \epsilon) > 1 - \frac{\text{Var}X}{\epsilon^2}$$

bzw.

$$P(|X - EX| < \epsilon) > 1 - \frac{\text{Var}X}{\epsilon^2}.$$

Also wäre der Anteil der korrekt befüllten Tassen mindestens 8/9.

2.4.

Mehrdimensionale Zufallsvariablen

<< hoch zum Anfang vom Inhaltsverzeichnis

<< hoch zum Kapitel 2 vom Inhaltsverzeichnis

Beispiel für mehrdimensionale Zufallsvariablen

In einer Studie über "Total Quality Management" wurde eine umfangreiche Befragung bei europäischen Produktionsbetrieben durchgeführt. Man erfasste zum einen den Aufwand für die Qualitätskontrolle während der laufenden Produktion in Prozent der

	y _j	0	5	10	15	P(X = X _j)
x _i						
0		0,00	0,00	0,10	0,30	0,4

Produktionskosten und des Weiteren die Aufwendungen für Reklamationen prozentual zum Umsatz.

Wir definieren die Zufallsvariablen:

X: Anteilige Kosten der Qualitätskontrolle [%]. Y: Anteilige Kosten der Reklamationen [%].

Es ergibt sich die unten dargestellte **gemeinsame**

Wahrscheinlichkeitstabelle mit der i-ten Zeile ($i = 1, \dots, n$) und der j-ten Spalte ($j = 1, \dots, m$). (Siehe darunter die graphische Darstellung der Tabelle.) Man erkennt, wie bei steigendem Aufwand der Qualitätskontrolle die Ausgaben für die Reklamationen sinken.

Die **gemeinsame Wahrscheinlichkeit** $P(X = 5 \wedge Y = 10) = 0,05$ werde bezeichnet als $f_{X,Y}(5;10)$.

Die spalten- bzw. zeilenweisen Summen der gemeinsamen Wahrscheinlichkeiten ergeben die **Randwahrscheinlichkeiten** oder auch **Einzelwahrscheinlichkeiten** der Zufallsvariablen X bzw. Y.

Es ergeben sich also für diese beiden Variablen die Wahrscheinlichkeitsverteilungen

x_i	0%	5%	10%
$f_X(x_i)$	0,4	0,2	0,4

y_j	0%	5%	10%	15%
$f_Y(y_j)$	0,2	0,2	0,2	0,4

Die Einzelwahrscheinlichkeit berechnet sich als

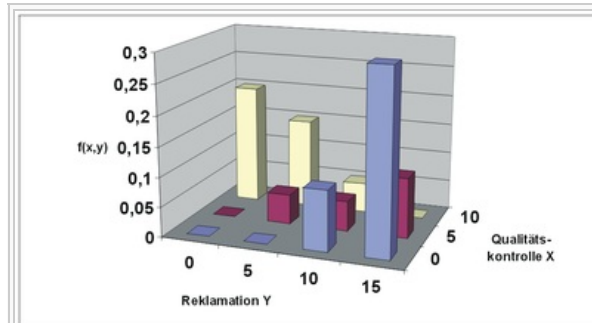
$$P(X = x_1) = f_X(x_1) = \sum_{j=1}^m f_{X,Y}(x_1; y_j) ,$$

also hier

$$P(X = 0) = f_X(0) = 0 + 0 + 0,1 + 0,3 = 0,4 .$$

5	0,00	0,05	0,05	0,10	0,2
10	0,20	0,15	0,05	0,00	0,4
$P(Y = Y_j)$	0,2	0,2	0,2	0,4	1,0

Gemeinsame Wahrscheinlichkeit von Qualitätskontrolle X und Reklamationskosten Y



Gemeinsame Wahrscheinlichkeit von Qualitätskontrolle X und Reklamationskosten Y

2.4.1.

Abhängigkeit von Zufallsvariablen

<< hoch zum Anfang vom Inhaltsverzeichnis

<< hoch zum Kapitel 2 vom Inhaltsverzeichnis

Stochastische Unabhängigkeit

Falls X und Y stochastisch unabhängig sind, ist

$$f_{X,Y}(x_i; y_j) = f_X(x_i) \cdot f_Y(y_j).$$

Beispiel:

Z.B. ist $P(X = 0 \wedge Y = 0) = 0$, aber $P(X = 0) \cdot P(Y = 0) = 0,4 \cdot 0,2 \neq 0$.

Also sind X und Y stochastisch abhängig. Es genügt schon, wenn die Unabhängigkeitsvoraussetzung für ein Paar nicht erfüllt ist.

Kovarianz

Man interessiert sich bei gemeinsam verteilten Variablen im allgemeinen auch dafür, inwieweit zwischen diesen Variablen ein Zusammenhang besteht. In unserer Wahrscheinlichkeitstabelle des Beispiels "Qualitätskontrolle" stehen beispielsweise links unten und rechts oben die größeren Wahrscheinlichkeiten, also scheinen niedrige Ausprägungen von X eher mit hohen Ausprägungen von Y und hohe Ausprägungen von X eher mit niedrigen Ausprägungen von Y einherzugehen.

Ein Maß für einen linearen Zusammenhang zweier Zufallsvariablen X und Y ist beispielsweise die Kovarianz $\text{cov}XY$. Sie ist für diskrete Zufallsvariablen definiert als

Wahrscheinlichkeitstabelle des Beispiels von oben Gemeinsame Wahrscheinlichkeit von Qualitätskontrolle X und Reklamationskosten Y					
$x \backslash y$	0	5	10	15	$f_X(x)$
0	0,00	0,00	0,10	0,30	0,4
5	0,00	0,05	0,05	0,10	0,2
10	0,20	0,15	0,05	0,00	0,4
$f_Y(y)$	0,2	0,2	0,2	0,4	1,0

$$\text{cov}XY = \sum_{i=1}^n \sum_{j=1}^m (x_i - EX)(y_j - EY) f_{X,Y}(x_i; y_j)$$

bzw. wegen des Verschiebungssatzes

$$\text{cov}XY = \sum_{i=1}^n \sum_{j=1}^m x_i \cdot y_j \cdot f_{X,Y}(x_i; y_j) - EX \cdot EY$$

Es ergibt für unser Beispiel

$$EX = 0 \cdot 0,4 + 5 \cdot 0,2 + 10 \cdot 0,4 = 5$$

und

$$EY = 0 \cdot 0,2 + 5 \cdot 0,2 + 10 \cdot 0,2 + 15 \cdot 0,4 = 9$$

und damit die Kovarianz

$$\begin{aligned} \text{cov}XY &= (0-5)(0-9) \cdot 0 + (5-5)(0-9) \cdot 0 + (10-5)(0-9) \cdot 0,2 \\ &\quad + (0-5)(5-9) \cdot 0 + (5-5)(5-9) \cdot 0,05 + (10-5)(5-9) \cdot 0,15 \\ &\quad + (0-5)(10-9) \cdot 0,1 + (5-5)(10-9) \cdot 0,05 + (10-5)(10-9) \cdot 0,05 \\ &\quad + (0-5)(15-9) \cdot 0,3 + (5-5)(15-9) \cdot 0,1 + (10-5)(15-9) \cdot 0 \\ &= 0 + 0 + (-5) \cdot 0,1 + (-30) \cdot 0,3 + 0 + 0 + 0 + 0 \\ &\quad + (-45) \cdot 0,2 + (-20) \cdot 0,15 + 5 \cdot 0,05 + 0 = -21,25 \end{aligned}$$

Eine positive Kovarianz deutet daraufhin, dass eher ein proportionaler Zusammenhang zwischen X und Y besteht, eine negative Kovarianz dagegen, dass eher ein umgekehrt proportionaler Zusammenhang zwischen X und Y besteht.

Korrelationskoeffizient

Ist die Kovarianz null, sind die Zufallsvariablen unkorreliert, sonst korreliert.

Die Kovarianz ist nicht normiert. Ein normiertes Maß für den linearen Zusammenhang stellt der **Korrelationskoeffizient** nach BRAVAIS-PEARSON $\rho_{X,Y}$ dar, der definiert ist als

$$\rho_{X,Y} = \frac{\text{cov}XY}{\sqrt{\text{var}X} \sqrt{\text{var}Y}}.$$

Es gilt für den Korrelationskoeffizienten $\rho_{X,Y}$:

$$-1 \leq \rho_{X,Y} \leq 1.$$

Ist $\rho_{X,Y}$ 1 oder -1, besteht ein exakter linearer Zusammenhang zwischen X und Y.

Sind X und Y stochastisch unabhängig, ist $\text{cov}XY$ und damit $\rho_{X,Y}$ gleich null. Der Umkehrschluss ist nicht zulässig, da eine nichtlineare Abhängigkeitsstruktur zwischen X und Y bestehen kann, die vom Korrelationskoeffizienten nicht erfasst werden kann.

Beispiel:

Wir berechnen zunächst die Varianz von X als

$$\text{var}X = (0 - 5)^2 \cdot 0,4 + (5 - 5)^2 \cdot 0,2 + (10 - 5)^2 \cdot 0,4 = 20$$

und entsprechend die Varianz von Y als

$$\text{var}Y = 34.$$

Damit erhalten wir

$$\rho_{X,Y} = \frac{\text{cov}(X, Y)}{\sqrt{\text{var}X} \sqrt{\text{var}Y}} = \frac{-21,25}{\sqrt{20} \sqrt{34}} = -0,8149.$$

Bedingte Wahrscheinlichkeiten von Zufallsvariablen

Auch für Zufallsvariablen sind bedingte Wahrscheinlichkeiten angebar, nämlich

die bedingte Wahrscheinlichkeit einer Zufallsvariablen als

$$P(X \leq x_i | X \leq x_k) = \frac{P(X \leq x_i \wedge X \leq x_k)}{P(X \leq x_k)}$$

und die bedingte Wahrscheinlichkeit zweier Zufallsvariablen

$$P(X \leq x_i | Y \leq y_j) = \frac{P(X \leq x_i \wedge Y \leq y_j)}{P(Y \leq y_j)}.$$

Entsprechendes gilt für \geq und $=$.

Ebenso gilt:

Wenn X und Y stochastisch unabhängig sind, ist

$$P(X \leq x_i \wedge Y \leq y_j) = P(X \leq x_i) \cdot P(Y \leq y_j)$$

für **alle** i, j.

Beispiele:

$$P(Y \geq 15 | Y \geq 5) = \frac{P(Y \geq 15 \wedge Y \geq 5)}{P(Y \geq 5)} = \frac{P(Y \geq 15)}{P(Y \geq 5)} = \frac{0,4}{0,8} = 0,5$$

„Die Hälfte aller Unternehmen mit Reklamationskosten hatte mindestens 15% Aufwand.“

$$P(Y \geq 5 | X = 10) = \frac{P(Y \geq 5 \wedge X = 10)}{P(X = 10)} = \frac{0,15 + 0,05 + 0}{0,4} = 0,5$$

„Die Hälfte aller Unternehmen mit sehr viel Qualitätskontrolle hatte Reklamationskosten.“

Funktion einer Zufallsvariablen

Lineare Transformation einer Zufallsvariablen

Der Student Bert hat eine kleine schicke Apartmentwohnung, die er hin und wieder säubern muss. Die Intervalle der Reinigungsaktionen sind unterschiedlich und lassen sich folgendermaßen beschreiben: Die Zeit in Wochen, die nach der letzten Säuberungsaktion verstrichen ist, wird als Zufallsvariable X bezeichnet. Die Intervalle verteilen sich folgendermaßen:

Zahl der Wochen bis zur nächsten Putzaktion x_i	0	1	2	3	4	5
Wahrscheinlichkeit $f(x_i)$	0,1	0,2	0,2	0,3	0,1	0,1

X hat den Erwartungswert $EX = 2,4$ und die Varianz $2,04$. Rechnen Sie das zur Übung selber nach.

Wenn Bert putzen muss, hängt der Aufwand in Stunden von der Zahl der Wochen ab, die er seine Wohnung vernachlässigt hat. Er braucht jedesmal ca. 1 Stunde für das Bad und einmal Durchsaugen. Für die restlichen Arbeiten muss er pro verstrichener Woche noch eine halbe Stunde Arbeitszeit hinzugeben. Morgen kommen seine Eltern zu Besuch. Mit welcher Wahrscheinlichkeit muss Bert heute 2 Stunden putzen? Wie lange putzt er durchschnittlich jedes Mal?

Hier überlegen wir uns zunächst mal, dass die Putzzeit von der vorherigen „Karenzzeit“ X abhängt. Sie ist also auch eine Zufallsvariable. Man könnte sie so darstellen:

$$Y = 1 + 0,5 \cdot X$$

Wie ist nun Y verteilt? Y hängt direkt von X ab und wir erhalten die Wahrscheinlichkeitstabelle

Zahl der Wochen bis zur nächsten Putzaktion x_i	0	1	2	3	4	5
Aufgewendete Putzzeit y_i	1	1,5	2	2,5	3	3,5
Wahrscheinlichkeit $f(y_i)$	0,1	0,2	0,2	0,3	0,1	0,1

Man kann sofort sehen, dass Bernd mit einer Wahrscheinlichkeit von 20% 2 Stunden putzen wird.

Wir wollen nun Erwartungswert und Varianz von Y ermitteln. Der Erwartungswert berechnet sich wie gewohnt als

$$\begin{aligned} EY &= \sum_i y_i \cdot f(y_i) = 1 \cdot 0,1 + 1,5 \cdot 0,2 + 2 \cdot 0,2 \\ &\quad + 2,5 \cdot 0,3 + 3 \cdot 0,1 + 3,5 \cdot 0,1 \\ &= 0,1 + 0,3 + 0,4 + 0,75 + 0,3 + 0,35 = 2,2 \end{aligned}$$

Das bedeutet er putzt durchschnittlich 2,2 Stunden.

Die Varianz ergibt sich analog als

$$\begin{aligned} \text{var}Y &= \sum_i y_i^2 \cdot f(y_i) - (EY)^2 = 1^2 \cdot 0,1 + 1,5^2 \cdot 0,2 + 2^2 \cdot 0,2 \\ &\quad + 2,5^2 \cdot 0,3 + 3^2 \cdot 0,1 + 3,5^2 \cdot 0,1 - 2,2^2 \\ &= 0,1 + 0,45 + 0,8 + 1,875 + 0,9 + 1,225 - 2,2^2 = 0,51 \end{aligned}$$

Schön wäre es allerdings, wenn man die Parameter der Verteilung etwas einfacher ausrechnen könnte. Y hat die schöne Eigenschaft, dass es eine **lineare Transformation** von X ist der Art

$$Y = a + bX.$$

Bei linearen Transformationen wie oben gilt

$$EY = a + b \cdot EX$$

und

$$\text{var}Y = b^2 \cdot \text{var}X.$$

Rechnen wir nach:

$$EY = 1 + 0,5 \cdot EX = 1 + 0,5 \cdot 2,4 = 1 + 1,2 = 2,2$$

und

$$\text{var}Y = 0,5^2 \cdot \text{var}X = 0,25 \cdot 2,04 = 0,51.$$

Standardisierung

Eine spezielle lineare Transformation ist die Standardisierung einer Zufallsvariablen X durch

$$Z = \frac{X - EX}{\sqrt{\text{var}X}}.$$

Man kann nämlich Z so umformen:

$$Z = \frac{X}{\sqrt{\text{var}X}} - \frac{EX}{\sqrt{\text{var}X}} = a + bX$$

mit $b = \frac{1}{\sqrt{\text{var}X}}$ und $a = -\frac{EX}{\sqrt{\text{var}X}}$, denn Erwartungswert und Varianz von X sind Konstanten.

Es ist dann $EZ = 0$ und $\text{var}Z = 1$.

Nichtlineare Funktion einer Zufallsvariablen

Lakonisch könnte man sagen: Eine nichtlineare Funktion ist eine Funktion, die nicht linear ist. Man kann sie also nicht in der Form $Y = a + bx$ schreiben. Beispiele sind etwa

$$Y = X^2, \quad Y = \sin X, \quad Y = \sqrt{X}$$

Hier kann man die Parameter im Allgemeinen nur über die Verteilung der Zufallsvariablen bestimmen.

Beispiel

Es hat sich herausgestellt, dass der Aufwand an Putzmitteln (ml pro qm) in Abhängigkeit von der verstrichenen Zeit quadratisch steigt mit der Funktion

$$Y = 2 + 1 \cdot X^2$$

Zahl der Wochen bis zur nächsten Putzaktion x_i	0	1	2	3	4	5
Aufgewendete Putzmittel y_i	2	3	6	11	18	27
Wahrscheinlichkeit $f(y_i)$	0,1	0,2	0,2	0,3	0,1	0,1

Hier kann man Erwartungswert und Varianz von Y nur mit den bekannten Formeln ermitteln, etwa

$$\begin{aligned} EY &= \sum_i y_i \cdot f(y_i) = 2 \cdot 0,1 + 3 \cdot 0,2 + 6 \cdot 0,2 \\ &\quad + 11 \cdot 0,3 + 18 \cdot 0,1 + 27 \cdot 0,1 \\ &= 0,2 + 0,6 + 1,2 + 3,3 + 1,8 + 2,7 = 9,8 \end{aligned}$$

Lineare Funktionen mehrerer Zufallsvariablen

Zwei Variablen

Gegeben sind zwei Zufallsvariablen X_1 und X_2 mit den Verteilungsparametern EX_1 , $\text{var}X_1$ und EX_2 , $\text{var}X_2$. Außerdem sind die beiden Zufallsvariablen korreliert mit der Kovarianz $\text{cov}X_1X_2$. Es wird eine Zufallsvariable

$$Y = b_0 + b_1X_1 + b_2X_2$$

gebildet. Analog zu oben errechnet sich der Erwartungswert von Y durch

$$EY = b_0 + b_1 EX_1 + b_2 EX_2.$$

Die Varianz von Y setzt sich aus den Einzelvarianzen der Zufallsvariablen zusammen. Hinzu kommt noch die Kovarianz:

$$\text{var}Y = b_1^2 \text{var}X_1 + b_2^2 \text{var}X_2 + 2b_1b_2 \cdot \text{cov}X_1X_2.$$

Wenn die zwei Zufallsvariablen X_1 und X_2 stochastisch unabhängig sind, ist ihre Kovarianz Null. Dann reduziert sich die Formel für die Varianz auf

$$\text{var}Y = b_1^2 \text{var}X_1 + b_2^2 \text{var}X_2.$$

Beispiel

Die Versorgung mit Getränken in einem Fußballstadion mittlerer Größe wird bei Spielen von einem Gastronomieunternehmen betrieben. Man weiß aus Erfahrung, dass die Zahl der verkauften Bierbecher von der Zahl der vorbestellten Eintrittskarten abhängt, und zwar in unterschiedlicher Weise von einheimischen und auswärtigen Besuchern. Es sei X_1 : Zahl der bestellten Karten von Einheimischen und X_2 : Zahl der bestellten Karten von Auswärtigen.

Es hat sich herausgestellt, dass $EX_1 = 10.000, EX_2 = 1000$ und $\text{var}X_1 = 2000, \text{var}X_2 = 300$ sind.

Zudem sind X_1 und X_2 korreliert, denn je interessanter ein Spiel, desto mehr Einheimische und Auswärtige schauen das Spiel an. Es ist $\text{cov}X_1X_2 = 400$.

Die Zahl der verkauften Getränke lässt sich angeben als

$$Y = 100 + 2X_1 + 3X_2.$$

Es ist hier

$$EY = 100 + 2 \cdot 10000 + 3 \cdot 1000 = 23100$$

und

$$\text{var}Y = 2^2 \cdot 2000 + 3^2 \cdot 300 + 2 \cdot 2 \cdot 3 \cdot 400 = 15500$$

Mehr als zwei Variablen

Gegeben sind n Zufallsvariablen X_i ($i = 1, \dots, n$) mit den Erwartungswerten EX_i , den Varianzen $\text{var}X_i$ und den paarweisen Kovarianzen $\text{cov}X_1X_2, \text{cov}X_1X_3, \dots, \text{cov}X_{n-1}X_n$. $\text{cov}X_iX_j$ ($i < j; i = 1, \dots, n-1; j = i+1, \dots, n$). Es sei

$$Y = b_0 + b_1X_1 + b_2X_2 + \dots + b_nX_n = b_0 + \sum_{i=1}^n b_iX_i.$$

Dann erhalten wir für den Erwartungswert

$$EY = b_0 + b_1EX_1 + b_2EX_2 + \dots + b_nEX_n = b_0 + \sum_{i=1}^n b_iEX_i.$$

Die Varianz von Y können wir als Summe der Varianzen und paarweisen Kovarianzen ermitteln als

$$\text{var}Y = \sum_{i=1}^n b_i^2 \text{var}X_i + 2 \sum_{i=1}^{n-1} \sum_{j=i+1}^n b_i b_j \text{cov}X_iX_j.$$

und, falls die Zufallsvariablen X_i stochastisch unabhängig sind, als Varianz

$$\text{var}Y = \sum_{i=1}^n b_i^2 \text{var}X_i.$$

Bei den ausgewählten Verteilungen handelt es sich um theoretische Zufallsverteilungen. Das sind Verteilungen, deren Form durch eine allgemein bekannte Funktion beschrieben wird. Oftmals kann beobachtet werden, dass die Verteilung bestimmter Zufallsvariablen annähernd durch eine theoretische Verteilung dargestellt werden kann, z. B. das Gewicht von Hähnchen einer Geflügelzucht ist meistens annähernd normalverteilt. Meist haben diese Verteilungen bestimmte Vorzüge, sie können leicht berechnet werden, und man kann auch wahrscheinlichkeitstheoretische Folgerungen ziehen. Hier bekannt ist bereits die Dreiecksverteilung.

3.1.	Diskrete Verteilungen	<< hoch zum Anfang vom Inhaltsverzeichnis << hoch zum Kapitel 3 vom Inhaltsverzeichnis
3.1.1.	Binomialverteilung	<< hoch zum Anfang vom Inhaltsverzeichnis << hoch zum Kapitel 3 vom Inhaltsverzeichnis

Binomialverteilung

Das **Urnenmodell mit Zurücklegen** bestimmt die binomialverteilte Zufallsvariable.

Gegeben ist eine Urne mit zwei Sorten Kugeln. Man spricht von einer dichotomen (griech: zweigeteilten) Grundgesamtheit. Es sind insgesamt N Kugeln in der Urne und M Kugeln der ersten Sorte. Der Anteil der Kugeln erster Sorte ist also

$$\theta = \frac{M}{N},$$

($0 \leq \theta \leq 1$). Es werden n Kugeln mit Zurücklegen gezogen. Es ist die Zufallsvariable definiert:

X : Anzahl der Kugeln 1. Sorte unter den n gezogenen Kugeln.

Beispiele für binomialverteilte Zufallsvariablen

- In einer Urne befinden sich 3 schwarze und 12 weiße Kugeln. Es werden fünf Kugeln gezogen, wobei jede Kugel sofort wieder zurückgelegt wird (Modell mit Zurücklegen). Wir definieren X als Zahl der weißen Kugeln bei $n = 5$ Entnahmen.
- 10 mal Würfeln. X : Zahl der Würfe mit einer Augenzahl von mindestens 5.
- Einem sehr großen Fertigungslos von Kondensatoren werden 10 Kondensatoren entnommen. Erfahrungsgemäß sind 15% der Kondensatoren schadhaft. X : Zahl der schadhaften Kondensatoren.
- In einer Schulklasse mit 30 Schülern und Schülerinnen wird täglich ein Kind per Los zum Tafeldienst bestimmt. X : Zahl der Tage, die Paula innerhalb von $n = 40$ Tagen Tafeldienst machen musste.

Exkurs

Beispiel: Sie würfeln 5 mal. Mit welcher Wahrscheinlichkeit erhalten Sie zweimal Sechs?

Offensichtlich handelt es sich bei diesem Problem um ein Urnenmodell mit Zurücklegen. Es wäre beispielsweise die Wahrscheinlichkeit, dass die ersten zwei Würfe Sechs ergeben:

$$\theta = \frac{1}{6} \cdot \frac{1}{6} \cdot \frac{5}{6} \cdot \frac{5}{6} \cdot \frac{5}{6} = \left(\frac{1}{6}\right)^2 \cdot \left(\frac{5}{6}\right)^3 = 0,01608.$$

Es gibt aber noch mehr Möglichkeiten, zwei Sechsen zu erhalten, nämlich:

(FFFSS), (FFSFS), (FFSSF), (FSFFS), (FSFSF), (FSSFF), (SFFFS), (SFFSF) und (SFSFF).

Hier bedeuten S: eine Sechs wird gewürfelt, F: keine Sechs wird gewürfelt. Es gibt insgesamt

$$\binom{5}{2} = \frac{5 \cdot 4}{1 \cdot 2} = 10$$

verschiedene Möglichkeiten, zwei Sechsen zu erhalten. Wir erhalten für die gesamte Wahrscheinlichkeit $P(X = 2)$, dass bei fünf Versuchen genau zwei Sechsen resultieren:

$$P(X = 2) = \binom{5}{2} \left(\frac{1}{6}\right)^2 \cdot \left(\frac{5}{6}\right)^3 = 10 \cdot 0,01608 = 0,1608.$$

Formale Darstellung

Die Zufallsvariable X ist **binomialverteilt mit den Parametern n und θ** . Ihre Wahrscheinlichkeitsfunktion lautet ($0 \leq \theta \leq 1$)

$$P(X = x) = b(x|n; \theta) = \begin{cases} \binom{n}{x} \theta^x (1 - \theta)^{n-x} & \text{falls } x = 0, 1, \dots, n \\ 0 & \text{sonst.} \end{cases}$$

Der Binomialkoeffizient berechnet sich als

$$\binom{n}{k} = \frac{n!}{k!(n-k)!}, \quad k, n \in \mathbb{N}, \quad k, n \geq 0.$$

Siehe auch in der Wikipedia: [Binomialkoeffizient](#)

Die Verteilungsfunktion $P(X \leq a) = B(a|n; \theta)$ ergibt sich als Summe der Wahrscheinlichkeiten einer diskreten Zufallsvariablen, wie in [Zufallsvariablen](#) oder [Diskrete Zufallsvariablen](#) erläutert.

Wie man der obigen Formel entnehmen kann, ist zur Berechnung der Wahrscheinlichkeiten die Kenntnis von N und M nicht erforderlich, es genügt die Bekanntheit von θ .

Weitere Kennwerte der Binomialverteilung sind

$$EX = n \cdot \theta \quad \text{und} \quad \text{var}X = n \cdot \theta \cdot (1 - \theta).$$

Beispiel: Verkehrszählung

Der Anteil der LKWs an den Kraftfahrzeugen auf deutschen Autobahnen soll für unser Beispiel 20% betragen. Im Rahmen einer Verkehrszählung an einer Auffahrt der Autobahn werden während einer Stunde 5 einfahrende Fahrzeuge zufällig erfasst.

1. Mit welcher Wahrscheinlichkeit befinden sich 2 LKWs in einer Stichprobe?
2. In wieviel Prozent der Stichproben befanden sich mindestens 2 LKWs in einer Stichprobe?

Es handelt sich offensichtlich um ein Modell mit Zurücklegen, denn ein Fahrzeug kann theoretisch auch mehrmals diese Auffahrt nehmen. Da wir die Fahrzeuge in LKW und Nicht-LKW unterscheiden, ist die betrachtete Grundgesamtheit dichotom (zwei Sorten Kugeln in der Urne). Wir definieren als Zufallsvariable X : Zahl der LKWs bei fünf gezählten Fahrzeugen.

X ist also binomialverteilt mit den Parametern $n = 5$ und $\theta = 0,2$ (20%), in Kurzschreibweise

$$X \sim b(x | 5; 0,2).$$

Wir werden zunächst die Wahrscheinlichkeitsfunktion von X bestimmen:

$x = 0$	$\binom{5}{0} \cdot \left(\frac{1}{5}\right)^0 \cdot \left(1 - \frac{1}{5}\right)^{5-0} = 1 \cdot 1 \cdot \left(\frac{4}{5}\right)^5 = \frac{1024}{3125}$	0,32768
$x = 1$	$\binom{5}{1} \cdot \left(\frac{1}{5}\right)^1 \cdot \left(\frac{4}{5}\right)^4 = 5 \cdot \frac{4^4}{5^5} = 5 \cdot \frac{256}{3125}$	0,4096
$x = 2$	$\binom{5}{2} \cdot \left(\frac{1}{5}\right)^2 \cdot \left(\frac{4}{5}\right)^3 = 10 \cdot \frac{4^3}{5^5} = 10 \cdot \frac{64}{3125}$	0,2048
$x = 3$	$\binom{5}{3} \cdot \left(\frac{1}{5}\right)^3 \cdot \left(\frac{4}{5}\right)^2 = 10 \cdot \frac{4^2}{5^5} = 10 \cdot \frac{16}{3125}$	0,0512

$x = 4$	$\binom{5}{4} \cdot \left(\frac{1}{5}\right)^4 \cdot \left(\frac{4}{5}\right)^1 = 5 \cdot \frac{4}{5^5} = 5 \cdot \frac{4}{3125}$	0,0064
$x = 5$	$\binom{5}{5} \cdot \left(\frac{1}{5}\right)^5 \cdot \left(\frac{4}{5}\right)^0 = 1 \cdot \frac{1}{3125} \cdot 1$	0,00032

Wir erhalten dann die Wahrscheinlichkeitstabelle

x_i	0	1	2	3	4	5
$b(x_i; 5; 0,2)$	0,32768	0,4096	0,2048	0,0512	0,0064	0,00032

Wir können also die gesuchten Wahrscheinlichkeiten aus der Tabelle ablesen

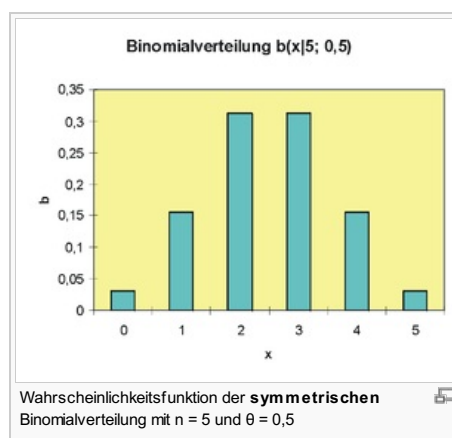
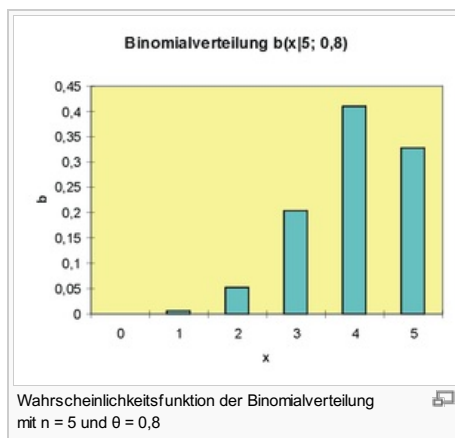
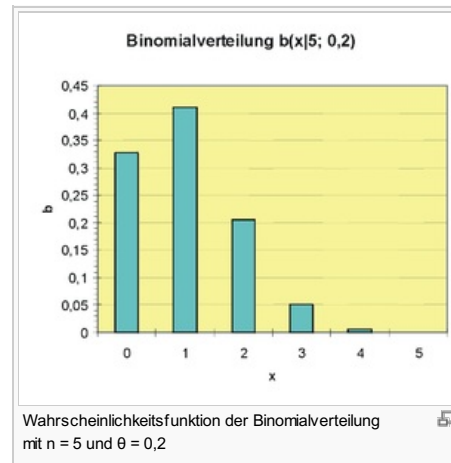
1. $P(X = 2) = 0,2048$
2. $P(X \geq 2) = 1 - P(X \leq 1) = 1 - (0,3277 + 0,4096) = 0,2627$

Eigenschaften der Binomialverteilung

Bei einem Urnenmodell mit Zurücklegen und zwei Sorten Kugeln (dichotome Grundgesamtheit) ist die Zahl der Kugeln erster Sorte bei n Entnahmen **immer** binomialverteilt.

Bei einem relativ kleinen Anteil θ ist die Verteilung rechtsschief (bzw. linkssteil), da die Wahrscheinlichkeit für ein kleines x groß ist. Bei einem relativ großen Anteil θ ist die Verteilung linksschief, da die Wahrscheinlichkeit für ein großes x eher groß ist.

Ist $\theta = 0,5$, ist die Verteilung symmetrisch bezüglich $x = \frac{n}{2}$.



Bemerkung

Bei großem n wird die Berechnung der Binomialkoeffizienten ein numerisches Problem, das allerdings beispielsweise mit der [Stirling-Formel](#) gelöst werden kann. Bei der Berechnung von Verteilungswerten kann allerdings die Addition der Wahrscheinlichkeiten sehr umständlich werden. Unter Umständen kann man die Funktionswerte der Binomialverteilung durch die [Poissonverteilung](#) oder auch durch die [Normalverteilung](#) approximieren.

Siehe auch in der Wikipedia: [Binomialverteilung](#)

3.1.2.

Hypergeometrische Verteilung

[<< hoch zum Anfang vom Inhaltsverzeichnis](#)

[<< hoch zum Kapitel 3 vom Inhaltsverzeichnis](#)

Das **Urnenmodell ohne Zurücklegen** bestimmt die hypergeometrisch verteilte Zufallsvariable.

Gegeben ist eine Urne mit zwei Sorten Kugeln. Man spricht von einer dichotomen (griech: zweigeteilten) Grundgesamtheit. Es sind insgesamt N Kugeln in der Urne und M Kugeln der ersten Sorte. Der Anteil der Kugeln erster Sorte ist also

$$\theta = \frac{M}{N},$$

($0 \leq \theta \leq 1$). Es werden n viele Kugeln ohne Zurücklegen gezogen. Es ist die Zufallsvariable definiert:

X : Anzahl der Kugeln 1. Sorte unter den n gezogenen Kugeln.

Beispiele für Hypergeometrische Verteilungen

- In einer Urne befinden sich 3 schwarze und 12 weiße Kugeln. Es werden fünf Kugeln ohne Zurücklegen gezogen (Modell ohne Zurücklegen). Wir definieren X als Zahl der weißen Kugeln bei $n = 5$ Entnahmen.
- Einem Fertigungslos von 100 Kondensatoren werden 10 Kondensatoren entnommen. Erfahrungsgemäß sind 15% der Kondensatoren schadhaft. X : Zahl der schadhaften Kondensatoren unter den 10 gezogenen.

Eine Zufallsvariable X ist **hypergeometrisch verteilt mit den Parametern N , M und n , wenn ihre Wahrscheinlichkeitsfunktion lautet**

$$P(X = x) = h(x|N; M; n) = \begin{cases} \frac{\binom{M}{x} \cdot \binom{N-M}{n-x}}{\binom{N}{n}} & \text{für } x = 0, 1, \dots, n \\ 0 & \text{sonst} \end{cases}$$

Die Verteilungsfunktion $P(X \leq a) = H(a|N; M; n)$ ergibt sich als Summe der Wahrscheinlichkeiten einer diskreten Zufallsvariablen, wie in [Zufallsvariablen](#) oder [Diskrete Zufallsvariablen](#) erläutert.

Weitere Kennwerte der hypergeometrischen Verteilung sind Erwartungswert und Varianz,

$$EX = n \cdot \frac{M}{N} \text{ und } \text{var} X = n \cdot \frac{M}{N} \cdot \left(1 - \frac{M}{N}\right) \frac{N-n}{N-1}.$$

Der letzte Bruch wird Korrekturfaktor genannt; er korrigiert die Varianz bei einem Modell ohne Zurücklegen. Wir können leicht sehen, dass für eine sehr große Grundgesamtheit (N) dieser Faktor etwa 1 wird. Bei einer großen Grundgesamtheit kann man also das Modell ohne Zurücklegen durch ein Modell mit Zurücklegen annähern.

Beispiel:

Von den sechs Rettichen, die eine Marktfrau auf dem Wochenmarkt verkauft, sind vier holzig. Der Student Paul sucht sich 4 Rettiche aus. Man könnte sich nun fragen: Mit welcher Wahrscheinlichkeit erwischt er alle holzigen?

Hier haben wir es unzweifelhaft mit einem **Modell ohne Zurücklegen** zu tun. Da wir holzige und nicht holzige Rettiche vor uns haben, ist die betrachtete Grundgesamtheit dichotom (zwei Sorten Kugeln in der Urne).

Wir definieren als Zufallsvariable X : Zahl der holzigen Rettiche bei $n = 4$ Entnahmen.

X ist also hypergeometrisch verteilt mit den Parametern $N = 6$, $M = 4$ und $n = 4$, in Kurzschreibweise

$$X \sim h(x|N; M; n) = h(x|6; 4; 4).$$

Wir werden zunächst die Wahrscheinlichkeitsfunktion von X bestimmen:

X = 0	$\frac{\binom{4}{0} \cdot \binom{6-4}{4-0}}{\binom{6}{4}} = \frac{1 \cdot 0}{15}$	0
X = 1	$\frac{\binom{4}{1} \cdot \binom{2}{3}}{\binom{6}{4}} = \frac{4 \cdot 0}{15}$	0
X = 2	$\frac{\binom{4}{2} \cdot \binom{2}{2}}{\binom{6}{4}} = \frac{6 \cdot 1}{15}$	$\frac{6}{15}$
X = 3	$\frac{\binom{4}{3} \cdot \binom{2}{1}}{\binom{6}{4}} = \frac{4 \cdot 2}{15}$	$\frac{8}{15}$
X = 4	$\frac{\binom{4}{4} \cdot \binom{2}{0}}{\binom{6}{4}} = \frac{1 \cdot 1}{15}$	$\frac{1}{15}$

Überlegen Sie sachlogisch, warum die ersten beiden Wahrscheinlichkeiten Null sind.

Der Student Paul wird also mit einer Wahrscheinlichkeit von 1/15 alle vier holzigen Rettiche erwischen.

Bemerkung

Werden M oder N groß, wird die Berechnung der Binomialkoeffizienten ein numerisches Problem, das allerdings beispielsweise mit der Stirling-Formel gelöst werden kann. Da der Unterschied zwischen einem Modell ohne Zurücklegen und mit Zurücklegen bei großem N unerheblich wird (ob man bei einer Entnahme 10000 oder 10001 Kugeln in der Urne hat, macht zahlenmäßig wenig aus), kann man bei großem N auch näherungsweise ein Modell mit Zurücklegen (siehe hierzu [Binomialverteilung](#)) verwenden. Häufig ist auch N unbekannt, hier kann das Modell ohne Zurücklegen gar nicht berechnet werden.

3.1.3.

Poissonverteilung

[<< hoch zum Anfang vom Inhaltsverzeichnis](#)

[<< hoch zum Kapitel 3 vom Inhaltsverzeichnis](#)

Wir betrachten eine poissonverteilte Zufallsvariable X mit den Ausprägungen 0, 1, 2,

Typische Beispiele für eine poissonverteilte Zufallsvariable sind:

- Es betreten in einer Minute durchschnittlich $\lambda = 2$ Kunden einen Kassenschalter. Wir definieren als X: Zahl der Kunden, die während einer bestimmten Minute an den Bankschalter kommen.
- Die Studentin Paula kauft sich in der Cafeteria ein Stück Rührkuchen. Wir definieren als X: Zahl der Rosinen in diesem Kuchenstück. Der Bäcker rechnet bei 20 Stück Kuchen mit 100 Rosinen. X ist also poissonverteilt mit dem Parameter $\lambda = 5$.
- Wir definieren als X: Zahl der Schadensfälle einer Versicherung im nächsten Jahr. Man weiß, daß pro Jahr durchschnittlich 500 000 Schadensfälle auftreten. Der Parameter ist hier $\lambda = 500\,000$.

Man geht also typischerweise von den folgenden Fragestellungen aus: Anzahl des Auftretens eines Phänomens in einer Zeit-, Gewichts- oder sonstigen Einheit. Die Zufallsvariable X ist **poissonverteilt** mit dem Parameter λ .

Ihre Wahrscheinlichkeitsfunktion lautet ($\lambda > 0$)

$$P(X = x) = p(x|\lambda) = \begin{cases} \frac{e^{-\lambda} \cdot \lambda^x}{x!} & \text{für } x = 0, 1, \dots \\ 0 & \text{sonst} \end{cases}$$

Die Verteilungsfunktion $P(X \leq a) = P_X(a|\lambda)$ ergibt sich als Summe der Wahrscheinlichkeiten einer diskreten Zufallsvariablen, wie in [Zufallsvariablen](#) oder [Diskrete Zufallsvariablen](#) erläutert.

Es gilt bei der Poissonverteilung: $EX = \text{var}X = \lambda$.

Die Poissonverteilung ist **reproduktiv**: Eine Summe von n stochastisch unabhängigen poissonverteilten Zufallsvariablen X_i ($i = 1, \dots, n$), mit jeweils dem Parameter λ_i , ist wiederum poissonverteilt, und zwar mit dem Parameter

$$\lambda = \sum_{i=1}^n \lambda_i$$

Beispiel:

Von den mundgeblasenen Gläsern einer Glashütte ist bekannt, dass im Durchschnitt 0,2 Fehler pro Glas auftreten.

Es ist die diskrete Zufallsvariable X : „Die Zahl der Unreinheiten in einem Glas“ annähernd poissonverteilt:

$$X \rightarrow p(x|0,2).$$

a) Mit welcher Wahrscheinlichkeit hat ein Glas genau einen Fehler?

$$P(X = 1) = \frac{e^{-0,2} \cdot 0,2^1}{1!} = 0,2 \cdot e^{-0,2} = 0,1637$$

b) Mit welcher Wahrscheinlichkeit hat ein Glas mindestens zwei Fehler?

$$P(X \geq 2) = 1 - P(X \leq 1) = 1 - \left(\frac{e^{-0,2} \cdot 0,2^0}{0!} + \frac{e^{-0,2} \cdot 0,2^1}{1!} \right)$$

$$= 1 - e^{-0,2} - 0,1637 = 1 - 0,8187 - 0,1637 = 0,0175.$$

c) Mit welcher Wahrscheinlichkeit enthalten drei Gläser zusammen mindestens zwei Fehler? Man geht davon aus, dass die Fehler der Gläser stochastisch unabhängig sind.

Man definiert als neue Zufallsvariable $Y = X_1 + X_2 + X_3$, mit X_1 als Zahl der Fehler des ersten Glases usw. Es ist dann $\lambda_y = 0,2 + 0,2 + 0,2 = 0,6$ und

$$P(Y \geq 2) = 1 - P(Y \leq 1) = 1 - \left(\frac{e^{-0,6} \cdot 0,6^0}{0!} + \frac{e^{-0,6} \cdot 0,6^1}{1!} \right)$$

$$= 1 - (e^{-0,6} + 0,6 \cdot e^{-0,6}) = 0,1219.$$

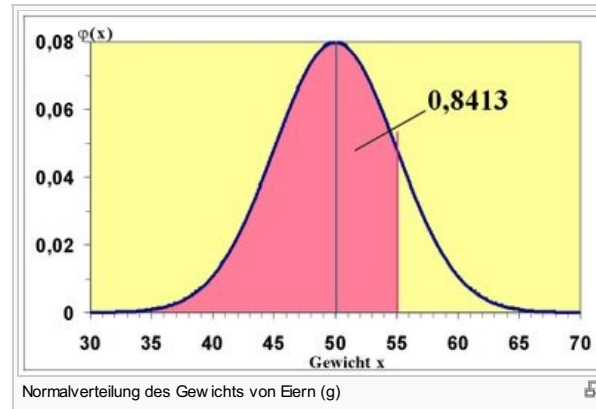
3.2.	Stetige Verteilungen	<< hoch zum Anfang vom Inhaltsverzeichnis << hoch zum Kapitel 3 vom Inhaltsverzeichnis
3.2.1.	Normalverteilung	<< hoch zum Anfang vom Inhaltsverzeichnis << hoch zum Kapitel 3 vom Inhaltsverzeichnis

Was ist die Normalverteilung?

Beispiel:

Auf einer Hühnerfarm mit sehr vielen Hühnern werden eine Woche lang die einzelnen Eier gewogen. Definieren wir die Zufallsvariable X : Gewicht eines Eis in Gramm. Es stellt sich heraus, daß ein Ei im Durchschnitt 50 g wiegt. Der Erwartungswert EX ist daher 50.

Außerdem sei bekannt, dass die Varianz $\text{var}X = 25 \text{ g}^2$ beträgt. Man kann die Verteilung des Gewichts annähernd wie in der Grafik darstellen. Man sieht, daß sich die meisten Eier in der Nähe des Erwartungswerts 50 befinden und daß die Wahrscheinlichkeit, sehr kleine oder sehr große Eier zu erhalten, sehr klein wird. Wir haben hier eine Normalverteilung vor uns. Sie ist typisch für Zufallsvariablen, die sich aus sehr vielen verschiedenen Einflüssen zusammensetzen, die man nicht mehr trennen kann, z.B. Gewicht des Huhns, Alter, Gesundheit, Standort, Vererbung usw.



Die Dichtefunktion der Normalverteilung ist definiert als

$$\phi_x(x|\mu;\sigma^2) = \frac{1}{\sqrt{2 \cdot \pi \cdot \sigma^2}} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right) \text{ für } x \in \mathbb{R},$$

wobei $E(X) = \mu$ und $\text{var}X = \sigma^2$ ist. Man sagt, X ist normalverteilt mit den Parametern μ und σ^2 , in Symbolschreibweise

$$X \sim \phi_x(x|\mu;\sigma^2) \text{ oder kürzer } X \sim N(\mu;\sigma^2).$$

In unserem Beispiel ist $X \sim N(50;25)$.

Die Normalverteilung ist symmetrisch bezüglich μ . Die Verteilung $P(X \leq a)$ von X ist wieder die Fläche unter dem Graph der Dichtefunktion. Sie wird bezeichnet als

$$P(X \leq a) = \Phi_x(a|\mu;\sigma^2) \text{ für alle } a \in \mathbb{R}.$$

Beispielsweise beträgt die Wahrscheinlichkeit, dass ein Ei höchstens 55 g wiegt, 0,8413. Das entspricht der roten Fläche in der Abbildung.

Das Integral der Dichtefunktion kann nicht analytisch berechnet werden. Die Werte der Verteilungsfunktion liegen i.a. [tabellarisch](#) vor. Es besteht nun das Problem, daß für jeden Wert von μ und σ^2 eine eigene Tabelle vorliegen müsste. Hier ist hilfreich, daß die aus X standardisierte Zufallsvariable Z wiederum normalverteilt ist und zwar mit den Parametern 0 und 1. Es kann jede beliebige Normalverteilung standardisiert werden. Mit Hilfe der standardisierten Zufallsvariablen wird dann die gesuchte Wahrscheinlichkeit bestimmt.

Standardnormalverteilung

Man definiert also eine neue Zufallsvariable

$$Z = \frac{X - EX}{\sqrt{\text{var}X}} = \frac{X - \mu}{\sigma}.$$

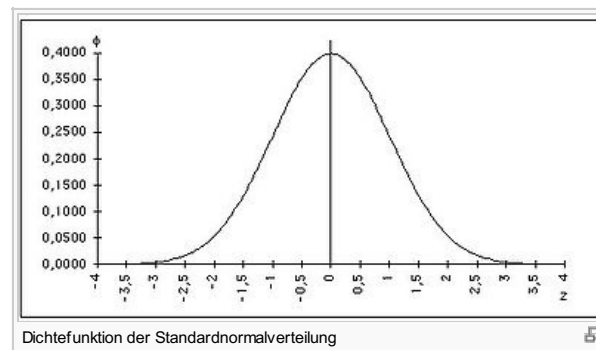
Diese Zufallsvariable Z ist normalverteilt mit $EZ = 0$ und $\text{var}Z = 1$. Ihre Dichtefunktion ist in der folgenden Grafik dargestellt. Es ist also $Z \sim N(0;1)$.

Die Dichtefunktion von Z ist

$$\phi_z(z|0;1) = \frac{1}{\sqrt{2 \cdot \pi}} \exp\left(-\frac{z^2}{2}\right) \text{ für } z \in \mathbb{R}$$

Ihre Verteilung, die man auch kurz als $\Phi(z)$ bezeichnet, ist (z const.)

$$P(Z \leq z) = \Phi(z) = \int_{-\infty}^z \phi(u) du.$$



Verteilungswerte

Es ist beispielsweise die Wahrscheinlichkeit

$$P(Z \leq 0,44) = 0,6700$$

und

$$P(Z > 1,64) = 0,0505.$$

Wir wollen nun den Anteil der Eier mit höchstens 55 g bestimmen, also $P(X \leq 55)$. Wir standardisieren:

$$z = \frac{x - \mu}{\sigma} = \frac{55 - 50}{\sqrt{25}} = 1.$$

Es ist dann

$$P(X \leq 55) = P(Z \leq 1) = \Phi(1) = 0,8413.$$

Der Wert 0,8413 der Verteilungsfunktion wird in der [Normalverteilungstabelle](#) ermittelt. Der folgende Ausschnitt aus der Tabelle soll die Vorgehensweise verdeutlichen: In der ersten Spalte der Tabelle sind die zwei ersten signifikanten Stellen der Ausprägung z angegeben, in der ersten Tabellenzeile die zweite Nachkommastelle, so dass sich beispielsweise $z = 1,00$ zusammensetzt aus $1,0 + 0,00$. Wo sich Zeile und Spalte des betreffenden Z-Wertes kreuzen, steht die gesuchte Wahrscheinlichkeit.

z	0,00	0,01	0,02
0,0	5000	5040	5080
0,1	5398	5438	5478
0,2	5793	5832	5871
0,3	6179	6217	6255
0,4	6554	6591	6628
0,5	6915	6950	6985
0,6	7257	7291	7324
0,7	7580	7611	7642
0,8	7881	7910	7939
0,9	8159	8186	8212
1,0	8413	8438	8461
1,1	8643	8665	8686
1,2	8849	8869	8888

Der errechnete Wert z kann gerundet werden, falls die errechneten Stellen die Zahl der Stellen des tabellierten z-Wertes übertreffen. Da die Verteilung von Z symmetrisch



bezüglich $\mu = 0$ ist, genügt die Tabellierung der Verteilungswerte ab $z = 0$ bzw. $\Phi(z) = 0,5$. Es gilt, wie man auch anhand der Grafik leicht sieht:

$$P(Z \leq -z) = P(Z \geq z) = 1 - P(Z \leq z) = 1 - \Phi(z)$$

bzw.

$$\Phi(-z) = 1 - \Phi(z).$$

Beispiel:

$$\begin{aligned} P(Z \leq -1) &= P(Z \geq 1) = 1 - P(Z \leq 1) \\ &= 1 - \Phi(1) = 1 - 0,8413 = 0,1587 \end{aligned}$$

Quantil

Häufig sucht man zu einer gegebenen Wahrscheinlichkeit p den dazugehörigen z -Wert $z(p)$. Er wird als **p-Quantil** bezeichnet.

Es gilt also:

$$P(Z \leq z(p)) = p.$$

Beispielsweise ist $z(0,975) = 1,96$. Es ist also hier die Wahrscheinlichkeit 0,975 **gegeben** und der dazugehörige z -Wert wird **gesucht**. Man sucht in der Tabelle die Wahrscheinlichkeit 0,9750 und bestimmt dann am Rand den betreffenden z -Wert 1,96.

Liegt p zwischen zwei Tabellenwerten, genügt es, als p den Tabellenwert zu verwenden, der p am nächsten liegt.

Beispiel:

Gesucht: $z(0,9)$

	näher bei	0,9	
Wahrscheinlichkeit Φ	0,8997		0,9015
z -Wert oder Quantil	1,28		1,29

Es ist also $z(0,9) \approx 1,28$.

Für eine Normalverteilung mit μ und σ^2 berechnet sich das p -Quantil als

$$x(p) = \mu + \sigma \cdot z(p).$$

Beispiel:

Wie schwer sind höchstens die 2/3 leichtesten Eier? Gesucht ist also $x(0,67)$:

$$x(p) = 50 + 5 \cdot z(0,67) = 50 + 5 \cdot 0,44 = 52,2.$$

Das schwerste der 67% leichtesten Eier wog also 52,2g.

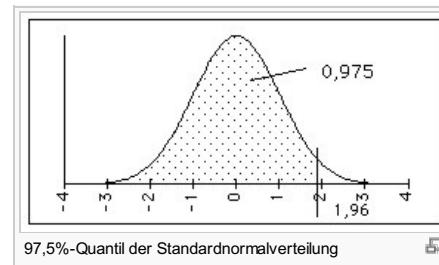
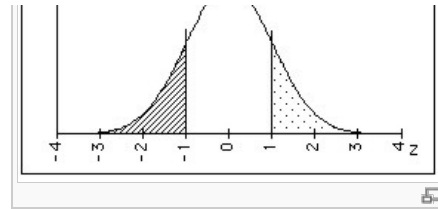
Übung zur Berechnung von $\Phi_Z(z)$

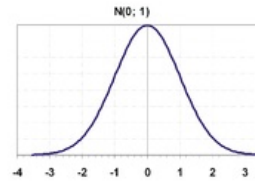
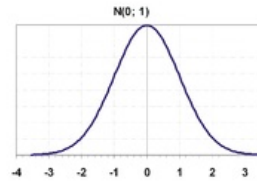
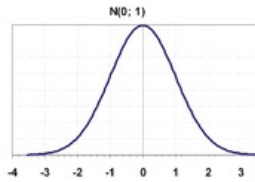
Schraffieren Sie die gesuchte Wahrscheinlichkeit in der Grafik und berechnen Sie die gesuchten Werte:

$$P(Z \leq 0,51)$$

$$P(Z \leq 2,0) =$$

$$P(Z \leq -0,51)$$

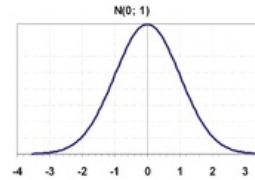
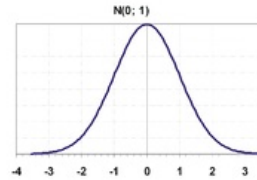
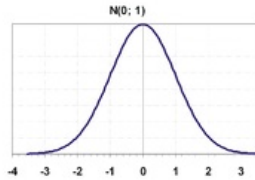




$$P(1,5 \leq Z \leq 2,35)$$

$$P(-0,8 \leq Z \leq 1,05)$$

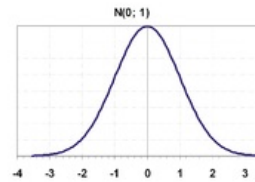
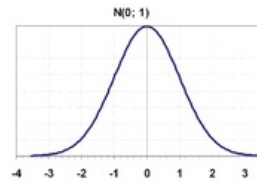
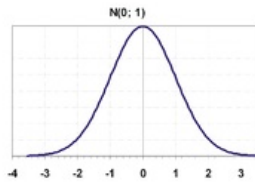
$$P(Z \geq -0,89)$$



$$P(Z \leq -1,68 \cup Z \geq 2)$$

$$P(Z \leq -1,96 \cup Z \geq 1,96)$$

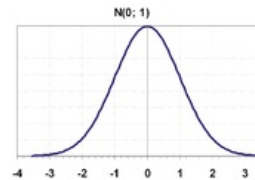
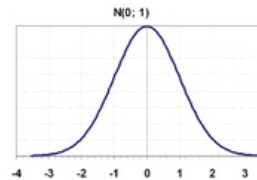
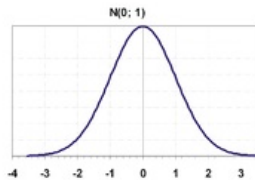
$$P(Z \leq -5)$$



$$z(0,975)$$

$$z(0,8)$$

$$z(0,2)$$



Übungen zum Eier-Beispiel

1. Wie groß ist die Wahrscheinlichkeit, daß ein Ei höchstens 60 g wiegt?
2. Wieviel Prozent der Eier wiegen höchstens 50 g?
3. Wie groß ist die Wahrscheinlichkeit, daß ein Ei mindestens 45 g wiegt?
4. Wieviel Prozent der Eier liegen zwischen 45 und 55 Gramm?
5. Mit welcher Wahrscheinlichkeit wiegt ein Ei genau 53 Gramm?
6. Welches Mindestgewicht haben die 30% schwersten Eier?

Lösungen:

Übung zur Berechnung von $\Phi_Z(z)$

a) 0,6950 b) 0,9772 c) 0,3050 d) 0,0574 e) 0,6412 f) 0,8133 g) 0,0693 h) 0,05 i) 0 j) 1,96 k) 0,84 l) -0,84

3.2.1.1.

Zentraler Grenzwertsatz

[<< hoch zum Anfang vom Inhaltsverzeichnis](#)

[<< hoch zum Kapitel 3 vom Inhaltsverzeichnis](#)

Gegeben sind die stochastisch unabhängigen Zufallsvariablen X_i ($i = 1, 2, \dots$).

Die Verteilungen der Summen Y_i

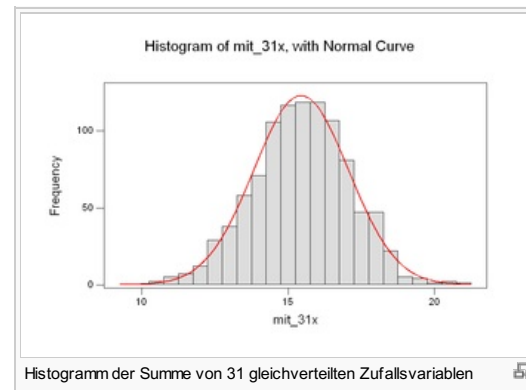
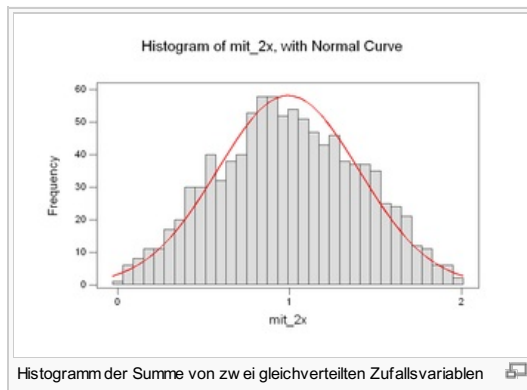
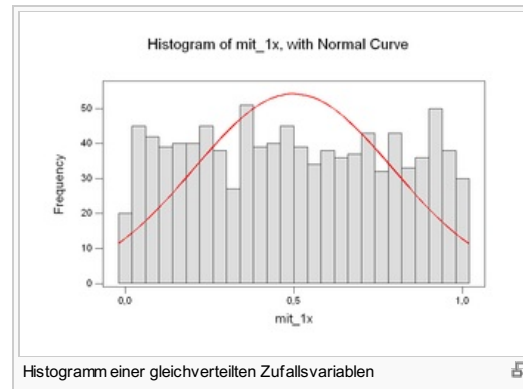
$$Y_1 = X_1, Y_2 = X_1 + X_2, \dots, Y_n = X_1 + X_2 + \dots + X_n, \dots$$

streben mit wachsendem n gegen die Normalverteilung. Als Faustregel gilt, daß die Verteilung einer Summe von mehr als 30 stochastisch unabhängigen Zufallsvariablen schon sehr gut annähernd mit der Normalverteilung bestimmt werden kann ($n > 30$).

Diese Regel ermöglicht zum einen die Bestimmung von Wahrscheinlichkeiten unbekannt verteilter Zufallsvariablen, zum anderen kann die Bestimmung kompliziert zu berechnender Wahrscheinlichkeitswerte mit der Normalverteilung angenähert (approximiert) werden.

Als Beispiel wurden je 1000 Zufallszahlen von im Intervall $[0;1]$ gleichverteilten Zufallsvariablen erzeugt. Der Graph ihrer Dichtefunktion bildet ein Rechteck.

Das Histogramm der Zufallszahlen lässt bei 1000 Werten deutlich das Rechteck erkennen. Bei der Summe von zwei gleichverteilten Variablen zeichnet sich die unimodale symmetrische Struktur schon deutlich ab, wobei zu bemerken ist, dass die Summe von zwei gleichverteilten Zufallsvariablen eine Dreiecksverteilung ergibt. Bei 31 Variablen ist die Näherung zur Normalverteilung schon sehr ausgeprägt.



3.2.1.2.

Lineare Funktionen der Normalverteilung

[<< hoch zum Anfang vom Inhaltsverzeichnis](#)

[<< hoch zum Kapitel 3 vom Inhaltsverzeichnis](#)

Linearkombinationen normalverteilter Zufallsvariablen

Gegeben sind n normalverteilte Zufallsvariablen X_i ($i = 1, \dots, n$), mit $X_i \sim N(\mu_i; \sigma_i^2)$. Die Linearkombination (lineare Funktion)

$$Y = a_0 + a_1 X_1 + a_2 X_2 + \dots + a_n X_n = a_0 + \sum_{i=1}^n a_i X_i$$

ist ebenfalls normalverteilt (Reproduktivität der Normalverteilung), und zwar mit dem Erwartungswert

$$EY = a_0 + \sum_{i=1}^n a_i EX_i = a_0 + \sum_{i=1}^n a_i \mu_i$$

und, falls die X_i ($i = 1, \dots, n$) stochastisch unabhängig sind, mit der Varianz

$$\text{var}Y = \sum_{i=1}^n a_i^2 \text{var}X_i = \sum_{i=1}^n a_i^2 \sigma_i^2.$$

Da die Varianz jedoch echt größer Null sein muss, muss zudem $a_j \neq 0$ für mindestens ein $j \in \{1, \dots, n\}$ gefordert werden.

Verteilung des Stichprobendurchschnitts Sind speziell die n Zufallsvariablen X_i ($i = 1, \dots, n$) sämtlich normalverteilt mit gleichem μ und gleichem σ^2 , ist die Linearkombination \bar{X} mit $a_0 = 0$, $a_1 = a_2 = \dots = a_n = 1/n$, also

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$$

normalverteilt dem Erwartungswert

$$E\bar{X} = \frac{1}{n} \sum_{i=1}^n \mu = \mu$$

und, falls die X_i ($i = 1, \dots, n$) stochastisch unabhängig sind, mit der Varianz

$$\text{var}\bar{X} = \frac{1}{n^2} \sum_{i=1}^n \sigma^2 = \frac{\sigma^2}{n}.$$

Beispiel

Die Firma Ziemlich&Unbekannt produziert die Güter I_x und $Ypsi$. Die monatliche Produktionsmenge schwankt zufällig, so dass für die produzierten Mengen die Zufallsvariablen definiert werden: X und Y [ME]. Man weiß:

$$X \sim N(20;5) \text{ und } Y \sim N(100;10).$$

Es wird vermutet, dass X und Y stochastisch unabhängig sind.

Wir interessieren uns für die monatlichen Gesamtkosten K in Crætos (C):

Die monatlichen Fixkosten betragen $a = 10.000$ C, die variablen Kosten für X : $b = 500$ C und für Y : $c = 200$ C.

Die monatlichen Gesamtkosten können also dargestellt werden als

$$K = a + bX + cY = 10000 + 500X + 200Y.$$

Wie ist also K verteilt? Wegen der Reproduktivitätseigenschaft der Normalverteilung müsste K wieder normalverteilt sein. Seine Parameter sind

$$EK = a + b EX + c EY = 10.000 + 500 \cdot 20 + 200 \cdot 100 = 40.000$$

und

$$\text{var}K = b^2 \text{var}X + c^2 \text{var}Y = 500^2 \cdot 5 + 200^2 \cdot 10 = 1.650.000.$$

Also ist $K \sim N(40.000; 1.650.000)$.

Mit welcher Wahrscheinlichkeit entstehen der Firma Gesamtkosten von mindestens 42.000 C?

Es ergibt sich

$$\begin{aligned} P(K \geq 42000) &= 1 - P(K \leq 42000) = 1 - \Phi_z\left(\frac{42000 - 40000}{\sqrt{1650000}}\right) \\ &= 1 - \Phi_z(1,57) = 1 - 0,9418 = 0,0582. \end{aligned}$$

χ^2 -Verteilung**Beispiel**

Wir haben 3 normalverteilte, paarweise stochastisch unabhängige Zufallsvariablen X_1 , X_2 und X_3 gegeben mit den Erwartungswerten μ_1 , μ_2 , μ_3 und den Varianzen σ_1^2 , σ_2^2 , σ_3^2 . Wir standardisieren diese Variablen und erhalten 3 standardnormalverteilte Zufallsvariablen Z_1 , Z_2 und Z_3 ,

$$Z_1 = \frac{X_1 - \mu_1}{\sigma_1}, Z_2 = \frac{X_2 - \mu_2}{\sigma_2}, Z_3 = \frac{X_3 - \mu_3}{\sigma_3}.$$

Nun werden die standardnormalverteilten Zufallsvariablen quadriert und aufsummiert. Wir erhalten eine neue Zufallsvariable

$$Y = Z_1^2 + Z_2^2 + Z_3^2.$$

Y ist χ^2 -verteilt mit 3 Freiheitsgraden.

Allgemein

Es gilt: Die Summe von m quadrierten, stochastisch unabhängigen, standardnormalverteilten Zufallsvariablen ist χ^2 -verteilt mit m Freiheitsgraden.

Man sieht anhand der Grafik, dass sich die Dichtefunktion mit wachsenden Freiheitsgraden einer symmetrischen Kurve nähert.

Die Wahrscheinlichkeit wird bezeichnet als $P(Y \leq a) = F_Y(a/n)$. Das p-Quantil ist $\chi^2(p;n)$.

Die Verteilungsfunktion der χ^2 -Verteilung kann nicht analytisch ermittelt werden. Numerische Berechnungen können beispielsweise aus Tabellenwerken, etwa [Tabelle der \$\chi^2\$ -Verteilung](#) ersehen werden. Da Y für jeden Freiheitsgrad eine eigene Verteilung besitzt, sind in kleineren Tabellen wie oben nur Quantile nach Freiheitsgraden und ausgewählten Wahrscheinlichkeiten aufgeführt. Es ist z. B. das 95%-Quantil (Spalte) der χ^2 -Verteilung mit 3 Freiheitsgraden (Zeile)

$F_Y(0,95;3) = 7,81$. Das bedeutet, die Wahrscheinlichkeit $P(Y \leq 7,81) = 0,95$.

Gilt $n > 30$, ist

$$Z = \sqrt{2X} - \sqrt{2n - 1}$$

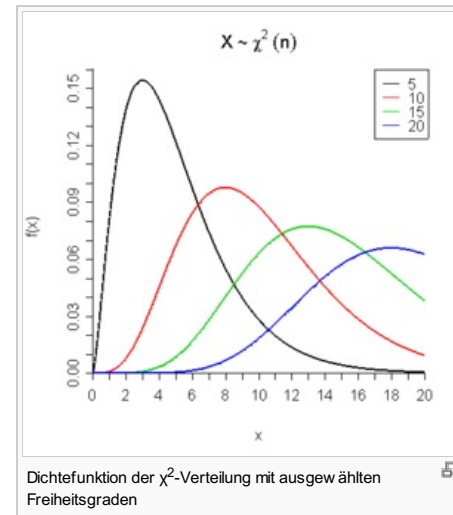
näherungsweise standardnormalverteilt.

Nähere Erläuterungen zur χ^2 -Verteilung, beispielsweise ihre Dichtefunktion, findet man bei [Wikipedia](#). Da die Dichtefunktion jedoch nicht für die Berechnung der Verteilungswerte unmittelbar verwendet werden kann, wird sie hier nicht angeführt.

Beispiele:

Sei Y χ^2 -verteilt mit 10 Freiheitsgraden. Es ist

- $P(Y \leq 15,99) = 0,9$
- $P(Y > 3,94) = 1 - P(Y \leq 3,94) = 1 - 0,05 = 0,95$
- $P(3,25 \leq Y \leq 20,48) = P(Y \leq 20,48) - P(Y \leq 3,25) = 0,975 - 0,025 = 0,95$



- 10%-Quantil von $Y : \chi^2(0,1;10) = 4,87$
- 95%-Quantil von $Y : \chi^2(0,95;10) = 18,31$

Sei Y χ^2 -verteilt mit 61 Freiheitsgraden. Gesucht ist $P(Y \leq 98)$. Hier ist die Zahl der Freiheitsgrade $k > 30$. Es wird eine neue Zufallsvariable $X = \sqrt{2Y}$ gebildet. X ist näherungsweise normalverteilt wie $N(\sqrt{2k-1}; 1) = N(11; 1)$. $P(Y \leq 98)$ entspricht also $P(X \leq \sqrt{2 \cdot 98}) = P(X \leq 14)$

Es ist $\Phi_X(14 | 11; 1) = \Phi_X\left(\frac{14 - 11}{1}\right) = \Phi_X(3) = 0,9987$.

Bemerkung

Die χ^2 -Verteilung ist **reproduktiv**, d. h. die Summe von zwei stochastisch unabhängigen χ^2 -verteilten Zufallsvariablen mit m und n Freiheitsgraden ist wieder χ^2 -verteilt mit $m+n$ Freiheitsgraden.

Die χ^2 -Verteilung ist eine so genannte Stichprobenverteilung.

Übung

1. Die Zufallsvariable X ist χ^2 -verteilt mit 12 Freiheitsgraden.
 1. Bestimmen Sie die Wahrscheinlichkeit, dass X kleiner als 6,30 ist.
 2. Bestimmen Sie die Wahrscheinlichkeit, dass X mindestens 18,55 beträgt.
 3. Bestimmen Sie das 5%-Quantil der Verteilung.
2. Die Zufallsvariable Y ist χ^2 -verteilt mit 40 Freiheitsgraden.
 1. Bestimmen Sie die Wahrscheinlichkeit, dass Y kleiner als 40 ist.
 2. Bestimmen Sie das 95%-Quantil der Verteilung.
3. Es sei $U = X + Y$.
 1. Bestimmen Sie den Erwartungswert von U .
 2. Bestimmen Sie die Wahrscheinlichkeit, dass U kleiner als 40 ist.

F-Verteilung

Beispiel

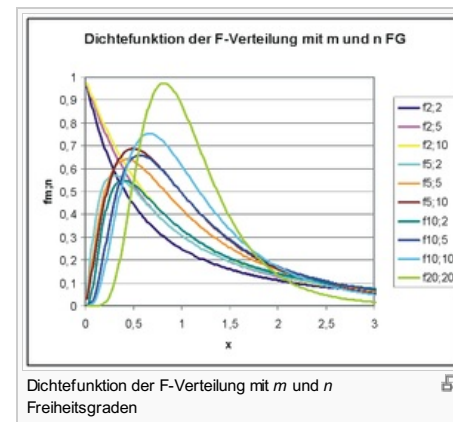
Wir haben die drei standardnormalverteilten Zufallsvariablen von oben und vier weitere Z_4, Z_5, Z_6 und Z_7 gegeben. Alle Variablen sind wieder stochastisch unabhängig. Der Quotient

$$F = \frac{\frac{Z_1^2 + Z_2^2 + Z_3^2}{3}}{\frac{Z_4^2 + Z_5^2 + Z_6^2 + Z_7^2}{4}}$$

ist dann F-verteilt mit 3 und 4 Freiheitsgraden.

Allgemein

Der Quotient aus zwei χ^2 -verteilten Zufallsvariablen, jeweils geteilt durch ihre Freiheitsgrade, wobei die Zufallsvariable im Zähler m und die im Nenner n Freiheitsgrade hat, ist F-verteilt mit m und n Freiheitsgraden. Einzelheiten dazu gibt es auch in der [Wikipedia](#). Man schreibt



$$F \sim F_{m;n}$$

Die Wahrscheinlichkeit wird bezeichnet als $P(F \leq a) = f_F(a|m;n)$. Das p-Quantil ist $F(p;m;n)$.

Auch die F-Verteilung liegt [tabelliert](#) vor und ist meistens nach ausgewählten Freiheitsgraden und Quantilen tabelliert. Eine nützliche Beziehung ist dabei

$$F(p;m;n) = \frac{1}{F(1-p;n;m)}.$$

Für viele Freiheitsgrade kann man sich die Faustregel merken: Sind m und n größer als 30, kann man die Quantile näherungsweise mit der Standardnormalverteilung ermitteln:

$$F(p;m;n) \approx z(p).$$

Die F-Verteilung ist ebenfalls eine Stichprobenverteilung. Sie ist aber nicht reproduktiv.

t-Verteilung Beispiel

Gegeben sind die standardnormalverteilten Zufallsvariablen von oben.

Der Quotient

$$t = \frac{Z_1}{\sqrt{\frac{Z_2^2 + Z_3^2 + Z_4^2 + Z_5^2}{4}}}$$

ist t-verteilt mit 4 Freiheitsgraden.

Allgemein

Der Quotient aus einer standardnormalverteilten Zufallsvariablen und der Wurzel einer χ^2 -verteilten Zufallsvariablen mit n Freiheitsgraden, geteilt durch ihre Freiheitsgrade, ist t-verteilt mit n Freiheitsgraden.

Die Wahrscheinlichkeit wird bezeichnet als $P(t \leq a) = f_t(a|n)$. Das p-Quantil ist $t(p;n)$.

Die Dichtefunktion der t-Verteilung ist, ähnlich wie die der Standardnormalverteilung, symmetrisch bezüglich des Erwartungswertes 0. Es gilt daher für die Berechnung der Verteilungswerte:

$$P(t \leq -a) = P(t \geq a),$$

mit

$$a \in \mathbf{R}.$$

Auch die t-Verteilung ist meistens nach Freiheitsgraden und ausgewählten Quantilen tabelliert: [t-Verteilung](#)

Für $n > 30$ kann man die Wahrscheinlichkeiten der t-Verteilung approximativ mit der Normalverteilung berechnen:

$$t(p;n) \approx z(p).$$

Bemerkungen:

- Das Quadrat einer t-verteilten Zufallsvariablen ist F-verteilt.
- Die t-Verteilung ist eine Stichprobenverteilung
- Weitere Eigenschaften können in der [Wikipedia](#) nachgelesen werden.

3.3.

Approximation von Verteilungen

[<< hoch zum Anfang vom Inhaltsverzeichnis](#)

[<< hoch zum Kapitel 3 vom Inhaltsverzeichnis](#)

Approximation: Approximation heißt **Näherung**, wie ja beispielsweise Alpha **Proxima** Centauri der uns am nächsten gelegene Stern ist. Wir wollen also Verteilungswerte, bei deren Berechnung wir heftige Unlustgefühle entwickeln, mit Hilfe anderer Verteilungen annähern. Sie werden nun mit Recht einwenden, dass das ja heutzutage mit der Entwicklung schneller Rechner eigentlich überflüssig sei. Nun hat man aber nicht immer einen Computer dabei (etwa in einer Klausur) oder es fehlt die Software zur Berechnung. MS-Excel bietet zwar solche Funktionen, aber

die Umsetzung ist etwas verquer, so dass häufig ein erhöhter Verstehensaufwand betrieben werden muss. Bei bestimmten Funktionswerten, wie großen Binomialkoeffizienten gehen schon mal Taschenrechner in die Knie.

Approximation diskreter Verteilungen durch diskrete Verteilungen

Die Wahrscheinlichkeitsfunktion der [Hypergeometrischen Verteilung](#) sieht so aus:

$$\frac{\binom{M}{x} \cdot \binom{N-M}{n-x}}{\binom{N}{n}}$$

Haben wir als Anwendung eine Kiste mit 10 Ü-Eiern gegeben, von denen 3 den gesuchten Obermotz enthalten, kann man etwa die Wahrscheinlichkeit, bei 5 Versuchen zwei Obermotze zu erhalten, leicht errechnen - naja, relativ leicht.

Aber betrachten wir den Fall: In einer Sendung von 500 speziellen Chips sind 100 Stück defekt. Bei der Eingangskontrolle werden 20 Chips getestet. Wenn jetzt die Wahrscheinlichkeit verlangt wird, dass genau 10 defekte Chips gezogen werden, erhält man

$$\frac{\binom{400}{10} \cdot \binom{100}{10}}{\binom{500}{20}}.$$

Spüren Sie schon Unlustgefühle? Vielleicht können wir uns hier die Berechnung mit der Binomialverteilung erleichtern, denn der Unterschied ist nur ob wir die getestete Chips zurücklegen oder nicht. Und bei 500 Chips spürt man kaum ob 20 davon heraus genommen sind ohne Zurücklegen, oder mit Zurücklegen. Vergleichen wir die beiden Verteilungen, fällt auf, dass beide den gleichen Erwartungswert haben: $EX = n\theta$. In den Varianzen aber unterscheiden sie sich,

Binomialverteilung:

$$\text{var} X = n\theta(1 - \theta)$$

und hypergeometrische Verteilung:

$$\text{var} X = n\theta(1 - \theta) \frac{N - n}{N - 1},$$

Der Unterschied liegt in dem Korrekturfaktor $(N-n)/(N-1)$. Wird nun N sehr groß, ist der Korrekturfaktor fast Eins und wir erhalten approximativ die Varianz der Binomialverteilung. Wie groß ist jetzt ein großes N ? Das kommt darauf an, wie genau wir die Näherung haben wollen. Für die Approximation der hypergeometrischen Verteilung durch die Binomialverteilung gibt es mehrere empfohlene Faustregeln, je nach Geschmack der Autoren. Eine der einfacheren Faustregeln, die man sich auch einigermaßen merken kann, ist

$$h(x|N; M; n) \approx b(x|n; \frac{M}{N}), \text{ wenn } \frac{n}{N} < 0,05$$

ist. Da in unserem Beispiel diese Voraussetzungen erfüllt sind, berechnen wir die gesuchte Wahrscheinlichkeit als

$$\binom{20}{10} \cdot 0,8^{10} \cdot 0,2^{10}.$$

Wir haben also das Modell ohne Zurücklegen durch ein Modell mit Zurücklegen angenähert. Man könnte so argumentieren: Wenn etwa 10000 Kugeln in einer Urne sind, macht es kaum einen Unterschied, ob beim 2. Versuch noch 9999 oder 10.000 Kugeln übrig sind. Analoges gilt für die Zahl der Kugeln 1. Sorte. Deshalb genügt auch die Angabe des Anteils θ dieser Kugeln an der Gesamtheit der Kugeln:

$$\theta = \frac{M}{N}.$$

Noch eine Bemerkung: Stellt man sich allerdings bei der Berechnung dieser Binomialkoeffizienten ein bisschen dumm an, protestiert die Software, weil man einen Überlauf erhält. Man kann allerdings hier mit der [Stirling-Formel](#) noch etwas ausrichten. Oder man [logarithmiert](#) die Fakultäten.

Für sehr kleines θ (oder sehr kleines $1-\theta$) und sehr großes n ist die Binomialverteilung wiederum annähernd Poisson-verteilt. Es ist nämlich die Poissonverteilung die Grenzverteilung der Binomialverteilung für $n \rightarrow \infty$ und $\theta \rightarrow 0$. Die Berechnung der Poissonverteilung ist einfacher als die Berechnung der Binomialverteilung. Eine Faustregel wäre hier etwa, dass eine binomialverteilte Zufallsvariable durch die Poisson-Verteilung angenähert werden kann, wenn $\theta \leq 0,05$ und $n \geq 50$ ist. Dann ist

$$b(x|N; M; n) \approx p(x|n\theta) .$$

Über den Umweg der Binomialverteilung kann dann auch die hypergeometrische Verteilung gegebenenfalls mit der Poisson-Verteilung approximiert werden:

$$h(x|N; M; n) \approx p(x|n\frac{M}{N}), \text{ wenn } \frac{n}{N} \leq 0,05, \theta \leq 0,05 \text{ und } n \geq 50$$

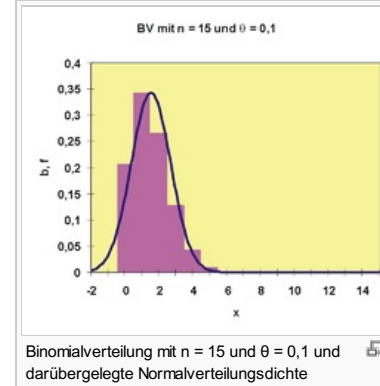
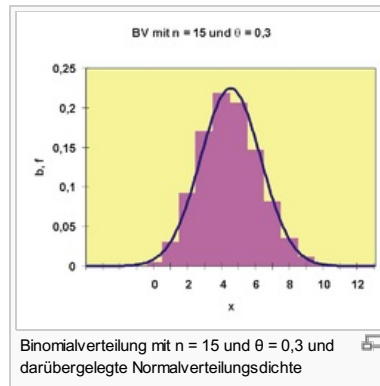
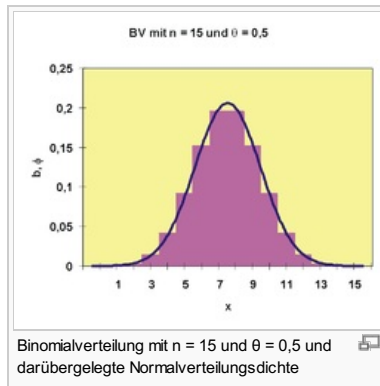
ist.

Weiter unten folgt eine tabellarische Zusammenfassung ausgewählter Approximationen.

Approximation diskreter Verteilungen durch die Normalverteilung

Was ist nun aber, wenn wir wissen wollen, wie groß die Wahrscheinlichkeit ist, dass höchstens 15 defekte Chips gefunden werden: $P(X \leq 15)$?

Hier müssen wir auf die oben beschriebene Weise 15 Wahrscheinlichkeiten ermitteln und addieren. Spätestens hier wünscht man sich eine Möglichkeit, so etwas schneller errechnen zu können. Es wäre doch angesagt, wenn man da die Normalverteilung verwenden könnte.



Vergleichen wir die Grafiken von den Binomialverteilungen. Es wurden hier die Wahrscheinlichkeiten als benachbarte Säulen dargestellt, was ja am optischen Erklärungswert nichts ändert.

Wir können deutlich erkennen, dass die Binomialverteilung für $\theta = 0,5$ symmetrisch ist. Hier passt sich die Normalverteilung am besten an. Je weiter θ von 0,5 abweicht, desto schlechter ist die Anpassung der Normalverteilung. Die so gut wie immer verwendete Faustregel ist, dass man mit der Normalverteilung approximieren darf, wenn

$$n > \frac{9}{\theta(1 - \theta)}$$

ist. Dürfen heißt natürlich nicht, dass es sonst polizeilich verboten ist, sondern dass sonst die Anpassung unbefriedigend ist.

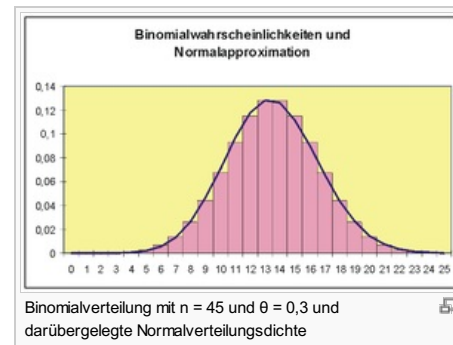
Eine Normalverteilung hat den Erwartungswert μ und die Varianz σ^2 . Wie soll man diese Parameter bei der Approximation ermitteln? Nun wissen wir ja, dass der Erwartungswert der Binomialverteilung und ihre Varianz

$$EX = n\theta \text{ und } \text{var}x = n\theta(1 - \theta)$$

sind, also nehmen wir doch einfach diese Parameter für die Normalverteilung, also

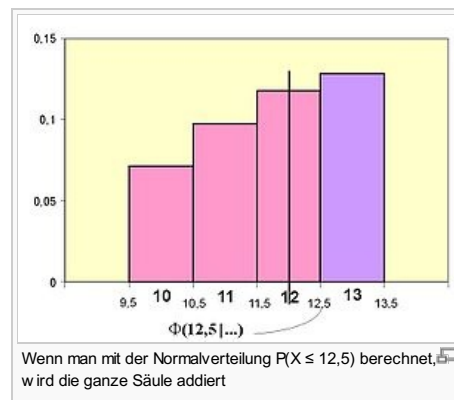
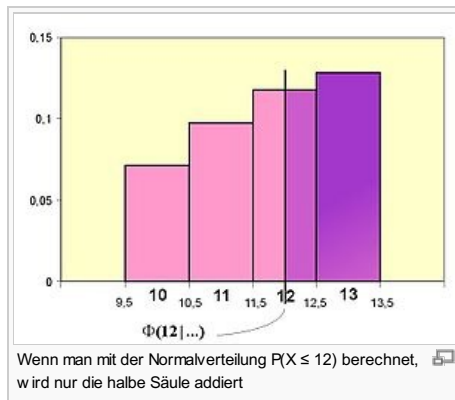
$$\mu = n\theta \text{ und } \sigma^2 = n\theta(1 - \theta).$$

Etwas fehlt uns noch: Wir nähern hier eine diskrete Verteilung durch eine stetige Verteilung an. Diskrete und stetige Verteilungen sind zwei völlig



unterschiedliche Konzepte. Wir betrachten hier das Beispiel einer Binomialverteilung mit $n = 45$ und $\theta = 0,3$.

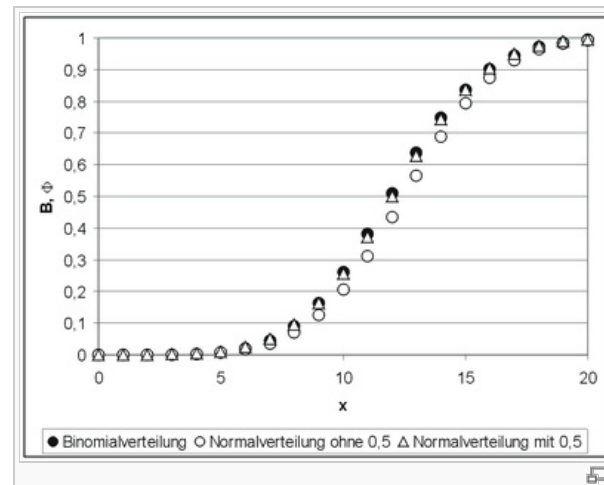
Nähern wir $P(X \leq 12) = B(12|45;0,3)$ durch $\Phi(12|45-0,3; 45-0,3-0,7)$ an, wird nur die halbe Säule addiert, denn die stetige Verteilung kennt keine Säulen. Soll die ganze Säule einbezogen werden, müssen wir bis 12,5 gehen, also $P(X \leq 12) = B(12|45;0,3)$ durch $\Phi(12,5|45-0,3; 45-0,3-0,7)$.



Den addierten Wert 0,5 nennt man **Stetigkeitskorrektur**.

Speziell gilt für die Wahrscheinlichkeit $P(X = a)$:

$$P(X = a) = b(a|n;\theta) \approx \Phi(a+0,5|n\theta; n\theta(1-\theta)) - \Phi(a-0,5|n\theta; n\theta(1-\theta)).$$



Approximation stetiger Verteilungen durch die Normalverteilung

Jetzt haben wir also auch noch stetige Funktionen, die wir mit der Normalverteilung annähern wollen. Was gibt es denn da für welche? Nun, welche die man oft braucht, etwa für [Schätzen und Testen](#), als da wären die χ^2 -Verteilung, die F-Verteilung und die t-Verteilung.

Nehmen wir uns doch mal die χ^2 -Verteilung vor. Ein Blick auf ihre [Dichtefunktion](#) verrät, dass diese mit wachsendem n immer symmetrischer wird, sich also der Normalverteilung annähert. Wir wissen, dass die χ^2 -Verteilung eine Summe von Zufallsvariablen, nämlich standardnormalverteilt, quadriert, ist und wir erinnern uns (gell?), dass nach dem zentralen Grenzwertsatz sich die Verteilung einer Summe von Zufallsvariablen der Normalverteilung annähert. Betrachten wir die mit n Freiheitsgraden χ^2 -verteilte Zufallsvariable X . Wir bilden eine neue Zufallsvariable

$$Y = \sqrt{2X}$$

Eine gängige Faustregel besagt für die Approximation für die Wahrscheinlichkeit $P(Y \leq y)$:

$$P(Y \leq y) \approx \Phi(y|\sqrt{2n-1}; 1) .$$

Die Dichtefunktion t-Verteilung dagegen hat eine ähnliche Form wie die Standardnormalverteilung, denn auch sie ist symmetrisch bezüglich der Null. Hier genügt eine einfache Faustregel: Wenn $n > 30$ ist, kann man die Verteilungswerte der t-Verteilung annähernd mit Hilfe der Standardnormalverteilung bestimmen:

$$t(x|n) \approx \Phi(x|0; 1) .$$

Tabelle der Approximationen

Gesuchte Verteilung	Approximation durch		
	Binomial	Poisson	Normal
$P(X \leq x)$ Binomial $B(x n\theta) \approx$	---	$P(x n\theta)$ falls $n \geq 50$ und $\theta \leq 0,05$	$\Phi(x + 0,5 n \cdot \theta; n \cdot \theta \cdot (1 - \theta))$ falls $n > \frac{9}{\theta(1 - \theta)}$
Hypergeometrische $H(x N; M; n) \approx$	$B(x n \frac{M}{N})$ falls $\frac{n}{N} < 0,05$	über Binomialverteilung	$\Phi(x + 0,5 n \cdot \frac{M}{N}; n \cdot \frac{M}{N} \cdot (1 - \frac{M}{N}) \cdot \frac{N - n}{N - 1})$ falls $n > \frac{9}{\frac{M}{N} \cdot (1 - \frac{M}{N})}$ und $\frac{n}{N} < 0,05$
Poisson $P(x \lambda) \approx$	---	---	$\Phi(x + 0,5 \lambda; \lambda)$ falls $\lambda > 9$
χ^2 -Verteilung $\chi^2(x n) \rightarrow$ $P(\sqrt{2X} \leq \sqrt{2x}) \approx$	---	---	$\Phi(\sqrt{2x} \sqrt{2n - 1}; 1)$ falls $n > 30$
t-Verteilung $t(x n) \approx$	---	---	$\Phi(x 0; 1)$ falls $n > 30$
F-Verteilung $F(x m; n) \approx$	---	---	$\Phi(x 0; 1)$ falls $m > 30$ und $n > 30$

4.

Deskriptive Statistik: Analyse eines Merkmals

[<< hoch zum Anfang vom Inhaltsverzeichnis](#)

[<< hoch zum Kapitel 4 vom Inhaltsverzeichnis](#)

Einführung

Die Verfahren der deskriptiven Statistik (beschreibende Statistik, empirische Statistik) haben als Grundlage die Erhebung bzw. Beobachtung von Daten. Es geht hier darum, diese Daten in geeigneter Weise zusammenzufassen, sie zu ordnen, sie grafisch darzustellen usw. Ziele der deskriptiven Statistik:

1. Die Daten einer empirischen Untersuchung möglichst übersichtlich zu präsentieren, so dass die wesentlichen Informationen schnell und optimal aufgenommen werden können. Beispiele: Tabellen, Säulendiagramme, Durchschnitte, Prognosen etc. Auf eine verteilungstheoretische Analyse wird verzichtet.
2. Man interessiert sich für die unbekannte Verteilung eines statistischen Merkmals, für Kennwerte der Verteilung usw. Da eine vollständige Erfassung dieses Merkmals meist zu teuer oder auch unmöglich ist, wird man sich auf eine Teilerhebung, eine Stichprobe, beschränken. Man schätzt nun mit Hilfe dieser Stichprobe die gesuchten Werte. Dabei versucht man, die Wahrscheinlichkeit einer Fehlschätzung miteinzubeziehen.

Analyse eines Merkmals

Die Analyse des Merkmals hängt u.a. davon ab, welche Informationen man wünscht:

- Verteilung: Ist sie symmetrisch oder schief, ein- oder mehrgipflig?
- Niveau der Daten, z.B. Durchschnitt, Median?
- Streuung der Einzelwerte: hoch oder niedrig?

- Sind mehrere Merkmale korreliert?

Definitionen in der deskriptiven Statistik

Beispiel:

Es wurden $n = 7$ Hunde befragt, wie gut ihnen das neue Fröhlix-Trockenfutter schmecke. Die Eingabe der Fragebögen in eine Datei ergab die unten folgende Liste. Anhand dieser Liste sollen Begriffe der deskriptiven Statistik erklärt werden.

Die Eigenschaften, die erhoben werden, sind die **Merkmale (statistische Variablen)** x, y, \dots . Das Objekt, dessen Eigenschaften erhoben (erfragt, gemessen) werden, ist die **Untersuchungseinheit (Merkmalsträger)**. Die Menge aller statistischen Einheiten ist die **Grundgesamtheit** (statistische Masse). Die möglichen Werte, die eine Variable annehmen kann, sind die **Ausprägungen (Realisationen)**. Die konkrete Merkmalsausprägung von x , die eine Untersuchungseinheit Nummer i aufweist, ist der **Merkmalswert (Beobachtungswert, Beobachtung)** x_i ($i=1,2, \dots, n$).

Name	Geschlecht <i>Merkmal</i> 1=w, 2=m u	Rasse x	Alter <i>Merkmal</i> y	Note für Futter 1, ..., 5 <i>Ausprägungen</i> z
Rex <i>Merkmalsträger</i>	2	Schäferhund	3	1
Rexona	1	Mischling	5	4 <i>Merkmalswert</i>
Lassie	1	Collie	1	2
Hasso	2	Neufundländer	2	1
Strolchi <i>Merkmalsträger</i>	2	Schnauzer	7	2
Susi	1	Spaniel	2	3
Waldi	2	Dackel	1 <i>Merkmalswert</i>	5

Es sind die Ausprägungen des Merkmals

Note: 1,2,3,4,5

und die Ausprägungen des Merkmals

Geschlecht: 1,2.

Skalierung des Merkmals

Beispiel

Grundlage des Beispiels ist die Hundetabelle von oben. Der Student Paul leistet beim Hersteller von Fröhlix ein Praktikum ab. Er soll die Ergebnisse der Befragung präsentieren. Er fasst die Hundetabelle von oben zusammen und erhält u.a.

Durchschnittliches Alter eines Hundes:

$$\frac{1}{7}(3 + 5 + 1 + 2 + 7 + 2 + 1) = \frac{21}{7} = 3.$$

Ein befragter Hund war also im Durchschnitt 3 Jahre alt.

Durchschnittliches Geschlecht eines Hundes:

$$\frac{1}{7}(2 + 1 + 1 + 2 + 2 + 1 + 2) = \frac{11}{7} = 1,57.$$

Ein Hund hatte also im Durchschnitt 1,57 Geschlecht. ????? Würden Sie den Studenten Paul später in dieser Firma einstellen?

Es ist natürlich höherer Schwachsinn, vom Merkmal Geschlecht den Durchschnitt zu bilden. Man kann damit keinen Durchschnitt bilden, denn seine Ausprägungen sind keine Zahlen. Geschlecht ist ein **qualitatives** Merkmal. Es ist anders **skaliert** als Alter.

Es gibt also Merkmale mit unterschiedlichen Messbarkeitsarten. Die Vorschrift für die Messbarkeit ist in einer **Skala** festgehalten.

Nominalskala

Merkmale wie

- Haarfarbe: braun, blond, ...;
- berufstätig ja/nein;
- Margarinemarke: Panorama, Botterama, ...

sind nominalsskaliert. Die Ausprägungen des nominalskalierten Merkmals können nicht **geordnet** werden, man kann sie nur **vergleichen** und **abzählen**. Es handelt sich um **qualitative** Merkmale. Erhalten die Ausprägungen Ziffern zugeordnet, handelt es sich nur um eine

Verschlüsselung (Codierung): 1 = männlich, 2 = weiblich.

Ordinalskala

Zwischen den Ausprägungen des ordinalskalierten (rangskalierten) Merkmals existiert eine Beziehung der Form mehr oder weniger, < oder >, besser oder schlechter o.ä., also eine Art **natürlicher Reihenfolge**.

Beispiele

- Sterne eines Hotels: *, **, ***; ...
- Beurteilung eines Produktes durch einen Konsumenten: Sehr gut, eher gut, eher schlecht, ganz schlecht
- Noten: 1, 2, 3, 4, 5

Für die Ausprägungen läßt sich also eine **Rangordnung** feststellen, aber die Abstände zwischen den Rängen sind nicht festgelegt. So ist beispielsweise die Note Vier nicht doppelt so schlecht wie Zwei.

Metrische Skala

Die Abstände zwischen den Ausprägungen des metrisch skalierten (quantitativen) Merkmals können **gemessen** werden. Es handelt sich bei den Ausprägungen um **(reelle) Zahlen**.

Beispiele: Kinderzahl, Einkommen, Temperatur, ...

Die metrischen Variablen werden noch in diskret und stetig unterschieden:

Ein Merkmal ist **diskret** (=unterschieden), wenn man die Ausprägungen abzählen kann.

Beispiel

- Kinderzahl: 0, 1, 2, ... , 20.
- Mein „Einkommen“, wenn ich falsch parke: 3 Euro (gesparte Parkgebühr) oder -10 Euro (Strafzettel).

Es gibt auch **abzählbar unendlich** viele Ausprägungen:

- Zahl der Ausfälle einer Maschine in einem Jahr: 0, 1, 2, ...

Ein Merkmal ist **stetig** (kontinuierlich), wenn sich in einem beschränkten Intervall der reellen Zahlen unendlich viele Ausprägungen (**überabzählbar viele**) befinden.

Beispiele: Wasserstand in einem Stausee; Gewicht eines Goldstücks; Temperatur; Körpergröße.

Bemerkung: Oft sind Merkmale eigentlich diskret, aber mit sehr vielen, nah beieinanderliegenden Ausprägungen, z.B. Einwohnerzahl, Preise (in Cents), usw. Hier definiert man das Merkmal zweckmäßigerweise als stetig, da es sich so besser analysieren lässt (**quasistetig**).

Übung

Wie sind die folgenden Merkmale skaliert?

- Täglicher Bierkonsum der Studentin Paula
 - - in Flaschen
 - - in Litern
- Bekenntnis: 1= röm.-kath., 2 = evang., 3 = sonst
- Gewicht von Bernhardinern
- Aufgabe: schwer - leicht
- Zahl der zustehenden Urlaubstage
- Jeansmarke

4.1.

Analyse von Daten eines metrischen Merkmals mit wenigen verschiedenen Beobachtungen

<< hoch zum Anfang vom Inhaltsverzeichnis
<< hoch zum Kapitel 4 vom Inhaltsverzeichnis

Behandelt wird hier ein metrisch skaliertes Merkmal, von dem nur wenige verschiedene Beobachtungen vorliegen, beispielsweise das Gewicht von 10 Schlachthähnchen oder die abgefüllte Menge von Kakao in 6 „250“-g Päckchen. Diese Konstellation wurde deshalb eigens hervorgehoben, weil sich damit viele Methoden der deskriptiven Statistik einfach erklären lassen.

4.1.1.

Häufigkeitsverteilung

<< hoch zum Anfang vom Inhaltsverzeichnis
<< hoch zum Kapitel 4 vom Inhaltsverzeichnis

Urliste

Beispiel

n = 10 „Pfundschalen“ Erdbeeren wurden nachgewogen. Es ergab sich durch Nachwiegen die **Urliste**

x ₁	x ₂	x ₃	x ₄	x ₅	x ₆	x ₇	x ₈	x ₉	x ₁₀
480	500	510	450	400	490	505	510	480	480

mit dem Merkmal x: Gewicht eines Schälchens (g). Die Werte wurden in der Reihenfolge der Erhebung, also ungeordnet, als Urliste erfasst. Diese Art der Darstellung ist unübersichtlich und daher nur für wenige Beobachtungen geeignet.

Urlisten können auch mehrere, unterschiedlich skalierte Merkmale enthalten. Beispielsweise ist die Tabelle mit den Hunden eine Urliste.

Häufigkeitsverteilung

Liegt ein metrisch skaliertes Merkmal oder ein ordinalskaliertes Merkmal mit vielen Ausprägungen vor, kann man zunächst einmal die Urliste der Größe nach ordnen, um einen gewissen Eindruck zu erhalten.

Beispiel

Die Indizes in den eckigen Klammern bedeuten, dass die Beobachtungen der Größe nach geordnet wurden.

x _[1]	x _[2]	x _[3]	x _[4]	x _[5]	x _[6]	x _[7]	x _[8]	x _[9]	x _[10]
400	450	480	480	480	490	500	505	510	510

Man erkennt nun, dass über die Hälfte der Schälchen untergewichtig waren.

Allerdings ist das Sortieren mühsam, fehleranfällig und doch nicht sehr informativ. Mit dem **Zweig-Blätter-Diagramm (stem-and-leaf display)** kann man jedoch sowohl metrische Beobachtungen relativ leicht sortieren als auch eine erste Häufigkeitsverteilung erzeugen.

Zweig-Blätter-Diagramm

Beispiel:

Für das Jahr 2003 liegt das reale Wachstum des Bruttoinlandsprodukts für 38 europäische Staaten vor (© Statistisches Bundesamt, Wiesbaden 200 [1] [↗](#))

4,7 1,1 3,9 -0,1 4,7 1,8 0,2 4,8 1,4 1,9 0,3 5,2 7,4 9,0 2,6 0,4 0,7 7,2 -0,8 0,3 0,7 3,7 -1,3 4,9 7,3 1,6 -0,5 4,0 4,2 2,3 2,4 2,9 5,8 4,8 2,9 2,1 4,7 2,0

Wir wollen die Daten ordnen und einen Eindruck von der Häufigkeitsverteilung gewinnen. Dazu werden wir jetzt ein Zweig-Blätter-Diagramm oder, für Anglophile, ein Stem-and-Leaf-Display erzeugen.

Zuerst wird der Zweig gemacht - aus den Einsern:

```
-1 |
-0 |
 0 |
 1 |
 2 |
 3 |
 4 |
 5 |
 6 |
 7 |
 8 |
 9 |
```

Der zweite Wert ist 1,1, das Blatt 1 wird an die 1 gehängt

```
-1 |
-0 |
 0 |
 1 | 1
 2 |
 3 |
 4 | 7
 5 |
 6 |
 7 |
 8 |
 9 |
```

Dann hängen wir die Blätter an den Zweig, und zwar, indem wir von links nach rechts durch die Daten wandern:

Der erste Wert ist 4,7. Das Blatt 7 wird an den Zweig 4 gehängt

```
-1 |
-0 |
 0 |
 1 |
 2 |
 3 |
 4 | 7
 5 |
 6 |
 7 |
 8 |
 9 |
```

Es folgen 3,9 -0,1 4,7 1,8 ...

```
-1 |
-0 | 1
 0 |
 1 | 18
 2 |
 3 | 9
 4 | 77
 5 |
 6 |
 7 |
 8 |
 9 |
```

Schließlich erhalten wir

```

-1 |      3
-0 |     185
 0 |    234737
 1 |    18496
 2 |   6349910
 3 |      97
 4 |  77890287
 5 |      28
 6 |
 7 |     423
 8 |
 9 |      0

```

Diese Prozedur war schnell erledigt. Wir bekommen schon einen guten Eindruck von der Verteilung der Beobachtungswerte. Kippen wir das Diagramm um 90°, erkennen wir eine Art Säulendiagramm. Außerdem können wir nun die Werte schnell der Größe nach sortieren. Wir erhalten nun unser Stengel-Blätter-Diagramm:

```

-1 |      3
-0 |     158
 0 |   233477
 1 |   14689
 2 |  0134699
 3 |      79
 4 | 02777889
 5 |      28
 6 |
 7 |     234
 8 |
 9 |      0

```

Für unsere Zwecke ist das Stem-and-Leaf-Display jetzt ausreichend. Ein Stem-and-Leaf-Display kann aber auch noch mehr Einzelheiten enthalten. Ist die Zahl der erhobenen Beobachtungen sehr groß, können die Werte in Klassen tabellarisch zusammengefaßt werden. Diese Art der Analyse erfolgt weiter unten.

Summenfunktion

Beispiel

Ein Autohaus hat von seinen $n = 10$ Filialen die Zahl der verkauften LKWs des letzten Jahres vorliegen. Es folgt die Urliste mit den x_i geordnet:

$x_{[1]}$	$x_{[2]}$	$x_{[3]}$	$x_{[4]}$	$x_{[5]}$	$x_{[6]}$	$x_{[7]}$	$x_{[8]}$	$x_{[9]}$	$x_{[10]}$
1	5	8	8	10	10	10	17	22	24

Wir wollen die absolute Summenfunktion $S(a)$ bestimmen. $S(a)$ gibt an, wieviel Beobachtungen $x_i \leq a$ sind:

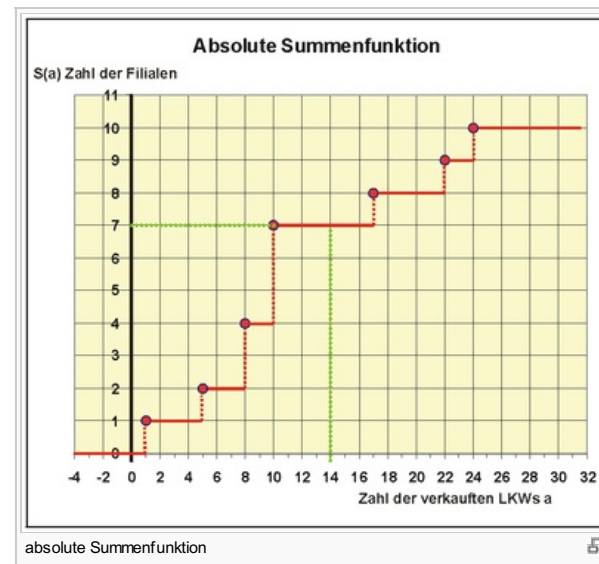
Zum Beispiel:

- $S(17) = 8$, also sind 8 Beobachtungen höchstens 17
- $S(8) = 4$, also gibt es 4 Filialen, die höchstens 8 LKWs verkauft haben

Wir leiten nun die Summenfunktion her, von links nach rechts:

- Zum Beispiel: $S(0,1) = 0$, denn keine Filiale hat höchstens 0,1 LKW verkauft. Ebenso ist $S(0,9) = 0$, usw... also

$$S(a) = 0 \text{ für } a < 1.$$



- Zum Beispiel: $S(1) = 1$, denn genau eine Filiale hat höchstens einen LKW verkauft. Ebenso ist $S(3) = 1$, denn es hat auch eine Filiale

höchstens drei LKWs verkauft. Ebenso $S(4,9999) = 1$..., also

$$S(a) = 1 \text{ für } 1 \leq a < 5.$$

- Zum Beispiel: $S(5) = 2$, also

$$S(a) = 2 \text{ für } 5 \leq a < 8.$$

usw... schließlich erhalten wir

$$S(a) = 10 \text{ für } a \geq 24.$$

Tragen wir die ermittelten Funktionswerte in die Grafik ein, sehen wir sofort, dass wir eine Treppenfunktion erhalten.

Die absolute Summenfunktion $S(a)$ ist die **Zahl** der Beobachtungen $x_i \leq a$. Die relative Summenfunktion gibt stattdessen die **Anteile** der Beobachtungen an der Urliste an:

$$S^*(a) = \frac{S(a)}{n}$$

Der Informationswert der kumulierten Häufigkeit $S(n)$ in der Grafik erschließt sich Ungeübten eher weniger. Aber man kann anhand der Grafik sofort Aussagen über die Verteilung machen. Man sieht beispielsweise sofort, daß z.B. 7 Werte kleiner als 14 sind, es haben also 70% der Filialen höchstens 14 LKWs verkauft.

4.1.2.

Lageparameter

[<< hoch zum Anfang vom Inhaltsverzeichnis](#)

[<< hoch zum Kapitel 4 vom Inhaltsverzeichnis](#)

Lageparameter

Der Lageparameter gibt an, auf welchem Niveau die Daten liegen.

Arithmetisches Mittel

Das arithmetische Mittel ist landläufig als „Durchschnitt“ bekannt. Es ist eigentlich nur für metrisch skalierte Merkmale (Problem Notendurchschnitt) geeignet. Es berechnet sich als

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

Beispiel Pfundschaalen Erdbeeren:

$$\begin{aligned} \bar{x} &= \frac{1}{10} (400 + 450 + 480 + 480 + 480 + 490 + 500 + 505 + 510 + 510) \\ &= \frac{4805}{10} = 480,5 \end{aligned}$$

Es waren die Schälchen also im Durchschnitt untergewichtig.

Median oder Zentralwert

Sind die Beobachtungswerte der Größe nach geordnet, also $x_{[1]}, x_{[2]}, x_{[3]}, \dots, x_{[n]}$, ist der Median z die Stelle, die die Teilgesamtheit in zwei gleiche Hälften teilt. Er kann für **rang-** und **metrisch** skalierte Merkmale verwendet werden.

n ungerade

Beispiel für $n = 7$

Es wurden 7 Autofahrer nach ihren Fahrtkosten befragt. Es ergab sich für das Merkmal x: Monatliche Ausgaben für Benzin (in Euro) die Liste

$x_{[1]}$	$x_{[2]}$	$x_{[3]}$	$x_{[4]}$	$x_{[5]}$	$x_{[6]}$	$x_{[7]}$
20	50	100	170	200	200	280

Es ist also der Median $z = 170$.

n gerade

Beispiel für $n = 10$ (Erdbeeren)

$x_{[1]}$	$x_{[2]}$	$x_{[3]}$	$x_{[4]}$	$x_{[5]}$		$x_{[6]}$	$x_{[7]}$	$x_{[8]}$	$x_{[9]}$	$x_{[10]}$
400	450	480	480	480	z	490	500	505	510	510

Der Median liegt zwischen dem 5. und 6. Beobachtungswert. Man nimmt hier den mittleren Wert

$$z = \frac{1}{2}(480 + 490) = 485.$$

Wir berechnen also den Median so:

n ungerade: z ist der $\frac{n+1}{2}$ te Wert $x_{[i]}$, also

$$z = x_{[\frac{n+1}{2}]}$$

n gerade: z liegt zwischen dem $\frac{n}{2}$ ten und dem $\frac{n}{2} + 1$ ten Beobachtungswert $x_{[i]}$, also

$$z = \frac{1}{2}(x_{[\frac{n}{2}]} + x_{[\frac{n}{2}+1]})$$

Bemerkungen:

- Der Median kann für ordinal- und metrisch skalierte Werte verwendet werden.
- Bei **sehr großem und geradem n** kann man vereinfachend

$$z = x_{[\frac{n}{2}]}$$

setzen.

Vergleich Median - arithmetisches Mittel

Beispiel:

Eine Autozeitschrift hat $n = 7$ PKWs einer bestimmten Marke getestet. Unter anderem wurde auch untersucht, ob das Auto zuverlässig anspringt.

Es ergab sich die geordnete Urliste

1 1 1 1 1 2 14

Wir erhalten als durchschnittliche Zahl der Startversuche

$$\bar{x} = \frac{1}{7}(1 + 1 + 1 + 1 + 1 + 2 + 14) = \frac{21}{7} = 3.$$

Wir würden hier also als Ergebnis erhalten: „Ein PKW sprang im Durchschnitt erst nach 3 Versuchen an“. Irgendwie erscheint einem das nicht gerechtfertigt. Bis auf einen PKW, der offensichtlich einen Ausreißer darstellt, sprangen ja alle Fahrzeuge zuverlässig an.

Wir verwenden nun den Median als Lageparameter: Der Median ist der 4. Wert, also $z = 1$. Hier ist also der Median eher zutreffend, doch so ganz zufrieden sind wir auch nicht, denn immerhin gab es ja auch 2 und 14 Versuche.

Wir sehen also, dass bei Verwendung des Median sehr viel Information der Daten verloren geht, andererseits reagiert aber das arithmetische Mittel empfindlich auf Ausreißer in den Daten.

Es gibt aber auch Kompromisse zwischen diesen beiden Extremen, beispielsweise das getrimmte Mittel:

$$x_T = \frac{1 + 1 + 1 + 1 + 2}{5} = \frac{6}{5} = 1,2$$

Es werden in der geordneten Urliste links und rechts jeweils ein oder mehrere Werte gestrichen. Aus den restlichen Beobachtungen berechnet man dann ein arithmetisches Mittel. Dieser Mittelwert erscheint eher die Sachlage zutreffend zu beschreiben. Man nennt Parameter, die nur schwach auf Ausreißer reagieren, resistente Parameter. Neben dem getrimmten Mittel gibt es noch mehrere andere Ansätze.

Der Vergleich des Medians mit dem arithmetischen Mittel kann als Ausreißeranalyse verwendet werden. Weicht der Median auffällig vom arithmetischen Mittel ab, sollten die Daten auf Ausreißer oder stark schiefe Verteilungen hin überprüft werden.

Weitere Lageparameter sind etwa der Modalwert, geometrisches Mittel oder harmonisches Mittel.

4.1.3.

Streuungsparameter

[<< hoch zum Anfang vom Inhaltsverzeichnis](#)

[<< hoch zum Kapitel 4 vom Inhaltsverzeichnis](#)

Varianz als Streuungsparameter

Der Lageparameter allein reicht für die Beschreibung einer Datenmenge nicht aus (analoges Problem wie bei Zufallsverteilungen). Information über die **Streuung** der Beobachtungswerte liefert ein **Streuungsparameter**. Es kommen verschiedene Kennwerte als Streuungsparameter in Betracht, beispielsweise die Varianz, die Spannweite, der Quartilsabstand und der Variationskoeffizient.

Varianz

Am häufigsten wird als Kennwert die Varianz verwendet, da sie wahrscheinlichkeitstheoretisch am besten zu untersuchen ist. Die Varianz sind die mittleren quadratischen Abweichungen der Einzelwerte x_i vom arithmetischen Mittel

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

Der Nenner $n-1$ wirkt vielleicht etwas befremdlich. Allerdings hat die Verwendung von $n-1$ statt n wahrscheinlichkeitstheoretische Vorzüge, wenn man die Varianz der Verteilung eines Merkmals mit s^2 schätzen möchte. Man nennt dieses Art der Varianz inferentielle Varianz.

Beispiel

Eine Firma möchte einen Kachelofen auf den Markt bringen, der für einen Komplettpreis zu erwerben ist. Für die Kalkulation dieses Preises benötigt die Firma Informationen über die Montagezeit für einen Kachelofen. Bei der Endmontage von 11 Kachelöfen ergaben sich die Zeiten

2,5 3 3 3,3 3,6 3 2,3 3 3,1 3,2 3

Die Varianz der Montagezeiten soll bestimmt werden. Nach der obigen Formel muss zunächst das arithmetische Mittel bestimmt werden:

$$\begin{aligned}\bar{x} &= \frac{1}{11}(2,5 + 3 + 3 + 3,3 + 3,6 + 3 + 2,3 + 3 + 3,1 + 3,2 + 3) \\ &= \frac{33}{11} = 3h\end{aligned}$$

Dann erhalten wir als Varianz

$$s^2 = \frac{1}{10}((2,5 - 3)^2 + (3 - 3)^2 + (3 - 3)^2 + \dots + (3 - 3)^2)$$

$$= \frac{1}{10}(0,25 + 0 + 0 + 0,09 + 0,36 + 0 + 0,49 + 0 + 0,01 + 0,04 + 0)$$

$$= \frac{1,24}{10} = 0,124h^2$$

Verzichtet man auf eine Schätzung, kann man auch die deskriptive Varianz

$$s_d^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

für die Beschreibung von statistischen Daten verwenden, was aber hier zur Vermeidung von Verwechslungen unterlassen wird.

Bei der manuellen Berechnung von s^2 ist es oftmals mühsam, erst die einzelnen Differenzen $x_i - \bar{x}$ zu bilden und dann zu quadrieren. Mit Hilfe des **Verschiebungssatzes** kann die laufende Differenzenbildung vermieden werden. Betrachten wir die Summe

$$\sum_{i=1}^n (x_i - \bar{x})^2.$$

Diese Summe lässt sich zerlegen in

$$Q = \sum_{i=1}^n x_i^2 - n \cdot \bar{x}^2.$$

Setzt man den Ausdruck oben ein, erhält man für die Varianz

$$s^2 = \frac{1}{n-1} Q = \frac{1}{n-1} \left(\sum_{i=1}^n x_i^2 - n \cdot \bar{x}^2 \right)$$

Beispiel:

$$s^2 = \frac{1}{10}(2,5^2 + 3^2 + 3^2 + \dots + 3^2 - 11 \cdot 3^2)$$

$$= \frac{1}{10} \cdot (100,24 - 99) = 0,124h^2$$

Da die Varianz ein quadratischer Ausdruck ist, hat sie z.B. auch die Einheit h^2 , wenn die x_i die Einheit h haben. Um die Varianz anschaulicher zu machen, kann man ihre Quadratwurzel, die Standardabweichung s betrachten:

Beispiel

$$s = \sqrt{0,124h^2} \approx 0,35h$$

also ca. 20 Minuten. Man könnte etwas flapsig sagen, dass die Montagezeit eines Ofens im Mittel 3 Stunden +/- 20 Minuten beträgt.

Auch die Varianz reagiert empfindlich auf Ausreißer. Es gibt hier resistente Streuungsparameter, die weiter unten behandelt werden.

4.2.

Analyse von Daten, die in Häufigkeitstabellen zusammengefasst sind

[<< hoch zum Anfang vom Inhaltsverzeichnis](#)

[<< hoch zum Kapitel 4 vom Inhaltsverzeichnis](#)

In den letzten Abschnitten lernten wir, wie man Daten eines stetigen, metrischen Merkmals, die als Urliste vorlagen, analysiert. Wir wollen nun Daten untersuchen, die man in Häufigkeitstabellen zusammenfassen kann. Im Gegensatz zur obigen Urliste können hier die Daten übersichtlich grafisch dargestellt werden. Man unterscheidet im Wesentlichen Daten eines metrischen Merkmals mit wenigen verschiedenen Ausprägungen und große Mengen von Daten mit vielen verschiedenen Ausprägungen, die man in Klassen zusammenfasst.

4.2.1.

Merkmale mit wenig Ausprägungen

[<< hoch zum Anfang vom Inhaltsverzeichnis](#)

[<< hoch zum Kapitel 4 vom Inhaltsverzeichnis](#)

Zu den **Merkmalen mit wenig verschiedenen Ausprägungen** gehören **nominal** skalierte, **ordinal** skalierte und **metrisch** skalierte Merkmale. Da sie nur **wenig Kategorien** haben, kann man sie in **Häufigkeitstabellen** zusammenfassen. Man nennt sie **häufbare Merkmale**.

Beispiele für Merkmale mit wenigen möglichen Ausprägungen:

- nominal skaliert: Augenfarbe von Studierenden
- ordinal skaliert: Note der Kundenzufriedenheit
- metrisch skaliert: Zahl der Autos in einem Haushalt

Bemerkung: Metrisch skalierte stetige Merkmale sind nicht unmittelbar häufbar, weil zu viele verschiedene Beobachtungen vorliegen.

4.2.1.1.

Verteilung

<< hoch zum Anfang vom Inhaltsverzeichnis

<< hoch zum Kapitel 4 vom Inhaltsverzeichnis

Wenn man vorliegende Daten analysiert, wird man sich zunächst für die Verteilung des Merkmals interessieren:

Ist die Verteilung der Variablen einigermaßen symmetrisch oder stark schief? Sind Ausreißer in den Daten? Ist die Verteilung eingipflig oder mehrgipflig? Der Statistiker freut sich meistens über eine symmetrische Verteilung, weil man hier häufig die Vorteile der Normalverteilung ausnützen kann.

Werkzeuge für die Analyse sind hierbei die Häufigkeitstabelle, die Summenfunktion und diverse Grafiken, denn bei einem Merkmal mit wenig Ausprägungen können attraktive Diagramme erstellt werden.

4.2.1.1.1.

Häufigkeit

<< hoch zum Anfang vom Inhaltsverzeichnis

<< hoch zum Kapitel 4 vom Inhaltsverzeichnis

Häufigkeitstabelle

Um eine Urliste von Beobachtungen eines Merkmals mit wenig Ausprägungen aufzubereiten, wird als erster Schritt der Analyse das Zählen des Auftretens der Realisationen stehen. Die Ergebnisse können in einer **Häufigkeitstabelle** zusammengefasst werden. Anhand der Daten eines nominalskalierten Beispiels wollen wir uns das Prinzip einer Häufigkeitstabelle ansehen.

Nominalskaliertes Merkmal

Beispiel

Es wurden 50 Personen telefonisch bezüglich gewisser Konsumpräferenzen befragt. Unter anderem erhob man den Familienstand. Es ist das Merkmal

x: Familienstand - mit den Ausprägungen 1=ledig, 2=verheiratet, 3=geschieden, 4=verwitwet.

Es ergab sich die Urliste

2 2 1 2 3 3 1 2 3 2 3 4 4 1 2 1 1 2 3 2 1 2 2 1 2
2 2 1 4 2 2 4 3 1 2 2 1 3 2 3 1 2 2 3 2 2 2 1 3 3

Wir wollen nun die Daten in einer Häufigkeitstabelle zusammenstellen:

j	Familienstand	absolute Häufigkeit	relative Häufigkeit
1	ledig	12	0,24
2	verheiratet	23	0,46
3	geschieden	11	0,22

4	verwitwet	4	0,08
Σ		50	1,00

Es sind insgesamt $n = 50$ Untersuchungseinheiten erhoben worden. Die (absoluten) Häufigkeiten n_j ($j = 1, \dots, 4$) verteilen sich auf $m = 4$ Kategorien (kategoriale Variable), wie in der Häufigkeitstabelle aufgelistet.

Wenn man sich für den Anteil der einzelnen Ausprägungen an der Gesamtheit interessiert, kann man auch die relativen Häufigkeiten bestimmen:

$$p_j = \frac{n_j}{n}$$

Es ist natürlich

$$\sum_{j=1}^m n_j = n \text{ bzw. } \sum_{j=1}^m p_j = 1$$

Für die Verteilung von Merkmalen mit wenig Ausprägungen kann man sehr ansprechende Grafiken erstellen.

Ordinalskaliertes Merkmal

Beispiel:

Bei der letzten Wiki-Matheklausur der Wikiversity ergaben sich die Noten wie folgt:

12 x 1, 15 x 2, 8 x 3, 3 x 4, 2 x 5

Hier erhält man die unten folgende Häufigkeitstabelle:

j	Note x_j	absolute Häufigkeit n_j	relative Häufigkeit p_j
1	sehr gut	12	12/40=0,3
2	gut	15	0,375
3	befriedigend	8	0,2
4	ausreichend	3	0,075
5	ungenügend	2	0,05
Σ		40	1

Auch hier bieten sich zur Veranschaulichung der Häufigkeiten Grafiken wie oben an.

Metrisch skaliertes Merkmal

Beispiel

Eine mainfränkische Weinbaustadt feiert ihr alljährliches Weinfest, bei dem auch die Winzerei Wavoma ihre Produkte anbietet. Sie verkauft Wein in Flaschen mit 0,5, 0,7, 1 und 2 Litern zu je 4, 5, 7 und 10 Euro. Es wurden am Sonntag Vormittag eingenommen (Merkmal x: Preis einer Flasche Wein (Euro)):

4 4 4 7 7 7 7 10 5 5 5 10 4 4 7 7 5 5 5 5 10 10 10 7

Wir erhalten die unten folgende Häufigkeitstabelle.

	Preis für eine Weinflasche	absolute Häufigkeit	relative Häufigkeit
--	----------------------------	---------------------	---------------------

j	x_j	n_j	p_j
1	4	5	$5/25=0,2$
2	5	8	0,32
3	7	7	0,28
4	10	5	0,2
Σ		25	1

Grafische Darstellungen

Eine weitere Art, Verteilungen eines Merkmals übersichtlich darzustellen, ist die grafische Darstellung. Mit hoher Aussagekraft der Grafik geht meist ein Informationsverlust einher, so daß die Grafik die Tabelle nicht ersetzen, sondern nur unterstützen kann.

Da Grafiken auf einen Blick informieren sollen, sollen sie nicht überladen sein. Häufig verwendet werden heute Piktogramme, d.h. Diagramme, die den Sachverhalt optisch anschaulich verdeutlichen.

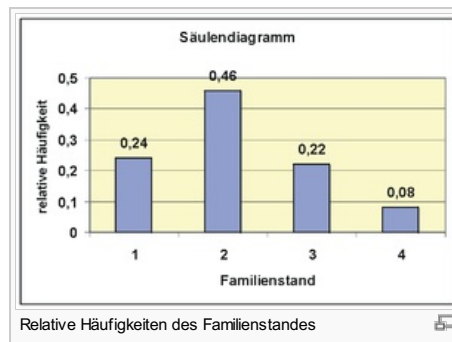
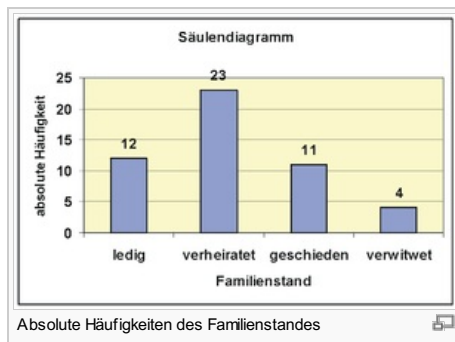
Für beliebig skalierte Merkmale mit wenigen Ausprägungen bieten sich eine Vielzahl grafischer Darstellungen an, darunter insbesondere Stabdiagramm, Säulendiagramm, Kreisdiagramm. Diese Diagramme eignen sich nicht für Urlisten mit vielen verschiedenen Beobachtungswerten.

Übung: Warum nicht?

Stabdiagramm bzw. Säulendiagramm

Auf der „x-Achse“ werden die verschiedenen Ausprägungen des Merkmals markiert. Dann werden die entsprechenden Häufigkeiten als Stab oder Säule senkrecht auf der Abszisse abgetragen.

Es sind hier anhand des obigen Beispiels bezüglich des Familienstandes die Säulendiagramme für die absoluten und relativen Häufigkeiten dargestellt. Wir sehen, dass die Struktur der Diagramme identisch ist.



Kreisdiagramm

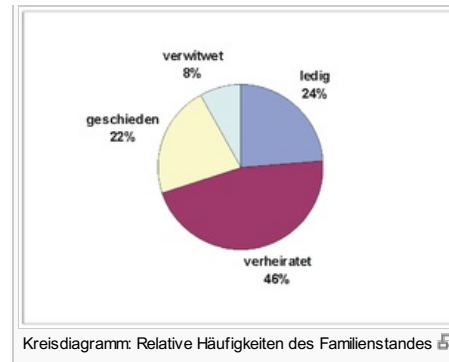
Im Kreisdiagramm wird die Gesamtfläche festgelegt. Die Häufigkeiten für die

Im Kreisdiagramm wird n als Gesamtfläche festgelegt. Die Häufigkeiten für die einzelnen Kategorien des Merkmals werden als „Tortenstücke“ eingetragen, wobei deren Fläche proportional zur Häufigkeit ist. Der zur Häufigkeit n_j gehörende Winkel α_j eines Segments berechnet sich dann aus der Verhältnisgleichung

$$\frac{\alpha_j}{360} = \frac{n_j}{n}$$

Sollen zwei verschiedene Gesamtheiten mit verschiedenen Gesamthäufigkeiten n_I und n_{II} mittels zweier Kreisdiagramme verglichen werden, kann man die Flächen der Kreise proportional zu den n_I und n_{II} darstellen.

Für die Darstellung von Kreisdiagrammen gibt es heutzutage genügend Anwendersoftware, so dass eine genauere Erläuterung unterbleiben kann.



4.2.1.1.2.

Summenfunktion

[<< hoch zum Anfang vom Inhaltsverzeichnis](#)
[<< hoch zum Kapitel 4 vom Inhaltsverzeichnis](#)

Summenfunktion

Man interessiert sich für Fragen wie „Wieviel % der Kunden gaben höchstens 5 Euro für eine Flasche Wein aus?“ oder „Wieviel Einwohner Deutschlands sind mindestens 65 Jahre alt?“. Man könnte nun die einzelnen Häufigkeiten einer Häufigkeitstabelle aufsummieren und so den Wert ermitteln, aber einfacher ist es, schon in der Häufigkeitstabelle die Häufigkeiten (abs. oder rel.) laufend aufzuaddieren. Es ergeben sich die **Summenhäufigkeiten** als **kumulierte Häufigkeiten** S_j (absolut) bzw. S_j^* (relativ). Aus den Summenhäufigkeiten lässt sich dann einfach die Summenfunktion bestimmen.

Summenhäufigkeiten sind nur sinnvoll, wenn man das Merkmal nach Größe ordnen kann, also nur bei ordinal oder metrisch skalierten Merkmalen. Aus der Summenhäufigkeit kann man die Summenfunktion ermitteln.

Beispiel der verkauften Weinflaschen

j	Preis für eine Weinflasche x_j	absolute Häufigkeit n_j	relative Häufigkeit p_j	absolute Summenhäufigkeit S_j	relative Summenhäufigkeit S_j^*
1	4	5	5/25=0,2	5	0,20
2	5	8	0,32	13	0,52
3	7	7	0,28	20	0,80
4	10	5	0,2	25	1,00
Σ		25	1		

Für die Erstellung der Summenfunktion müssen die Beobachtungen

Für die Erstellung der Summenfunktion müssen die Beobachtungen der Urliste geordnet vorliegen. Die Häufigkeitsverteilung enthält alle Werte der Urliste geordnet. Analog zu oben kann man sich beispielsweise überlegen:

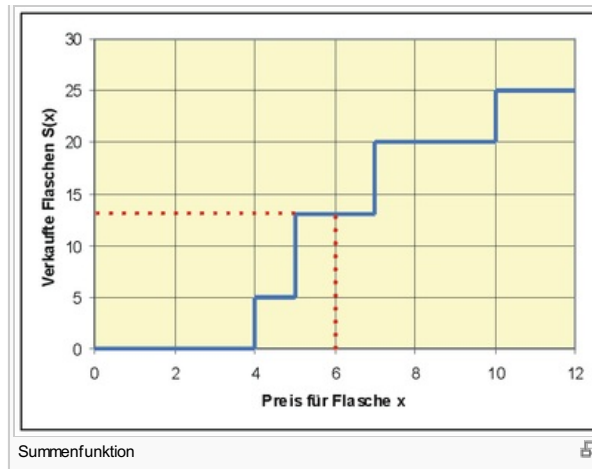
20 Kunden zahlten höchstens 7 Euro für eine Flasche, also $S(7) = 20$.

So können wir wieder wie vorher die Summenfunktion von links her aufbauen:

0 Kunden zahlten höchstens 2 Euro für eine Flasche, also $S(2) = 0$

usw.

Nun können wir die kumulierten Häufigkeiten auch aus der Grafik ablesen: z.B. $S(6) = 13$, es sind also 13 Flaschen zu einem Preis von höchstens 6 Euro verkauft worden.



4.2.1.2.

Lageparameter

<< hoch zum Anfang vom Inhaltsverzeichnis

<< hoch zum Kapitel 4 vom Inhaltsverzeichnis

Arithmetisches Mittel

Beispiel

Es wurden in einem Einkaufszentrum $n = 20$ Kunden bezüglich der Kinderzahl befragt. Wir erhielten die geordnete Urliste

0 0 0 0 1 1 1 1 1 2 2 2 2 2 3 3 3 4 5 5

Es resultierte die Häufigkeitsverteilung

j	Zahl der Kinder x_j	absolute Häufigkeit n_j	relative Häufigkeit p_j	$x_j n_j$	$x_j p_j$
1	0	4	0,2	0	0
2	1	5	0,25	5	0,25
3	2	5	0,25	10	0,5
4	3	3	0,15	9	0,45
5	4	1	0,05	4	0,2
6	5	2	0,1	10	0,5
Σ		20	1	38	1,9

Wir bestimmen das arithmetische Mittel als

$$\bar{x} = \frac{1}{20}(0 + 0 + 0 + 0 + 1 + 1 + 1 + 1 + 1 + 2 + 2 + 2 + 2 + 2 + 3 + 3 + 3 + 4 + 5 + 5) = \frac{38}{20} = 1,9$$

Wir können das Mittel aber auch so berechnen:

$$\bar{x} = \frac{1}{20}(4 \cdot 0 + 5 \cdot 1 + 5 \cdot 2 + 3 \cdot 3 + 1 \cdot 4 + 2 \cdot 5) = \frac{38}{20} = 1,9$$

was in Formelschreibweise ergibt

$$\bar{x} = \frac{1}{n} \sum_{j=1}^m x_j \cdot n_j.$$

Ermitteln wir das arithmetische Mittel von Hand, können wir in der Häufigkeitstabelle die Summanden $x_j n_j$ in der jten Zeile eintragen und aufsummieren.

Alternativ können wir das arithmetische Mittel mit Hilfe der relativen Häufigkeit p_j ermitteln:

$$\bar{x} = \sum_{j=1}^m x_j \cdot p_j.$$

Zur Verdeutlichung ist auch diese Variante in der Häufigkeitstabelle aufgeführt.

Für ordinal- oder nominalskalierte Merkmale ist das arithmetische Mittel nicht geeignet.

Entsprechende Überlegungen gelten auch für die Varianz s^2 der Stichprobe.

Median

Beispiel mit den verkauften Weinflaschen

Wir haben die Urliste nun **geordnet**.

4 4 4 4 4 5 5 5 5 5 5 5 5 7 7 7 7 7 7 7 10 10 10 10 10

Der Median teilt die kleineren 50% der Datenwerte von den 50% größeren Werten ab. Also liegt hier der Median auf dem 13. Beobachtungswert.

Bei Daten in Häufigkeitstabellen liegen die Werte schon zwangsläufig geordnet vor. Es muss nur die Kategorie gefunden werden, die den Median enthält.

j	Preis für eine Weinflasche x_j	absolute Häufigkeit n_j	absolute Summenhäufigkeit S_j
1	4	5	5
2	5	8	13
3	7	7	20
4	10	5	25
Σ		25	

Anhand der Summenhäufigkeiten können wir sehen, dass der 13. Wert gerade noch in der 2. Kategorie liegt. Diese Kategorie ist die **Einfallsklasse** des Medians.

4.2.1.3.

Streuungsparameter

<< hoch zum Anfang vom Inhaltsverzeichnis

<< hoch zum Kapitel 4 vom Inhaltsverzeichnis

Hier wollen wir die Berechnung der **Varianz eines häufbaren metrischen Merkmals** ansehen. Unsere Überlegungen laufen analog zum arithmetischen Mittel. Wir betrachten das

Beispiel mit den verkauften Weinflaschen

Aus der Urliste mit 25 Beobachtungen:

4 4 4 4 4 5 5 5 5 5 5 5 5 7 7 7 7 7 7 7 10 10 10 10 10

berechnen wir die Stichprobenvarianz aus

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2.$$

In dieser Formel ist x_i die i . Beobachtung aus der Urliste.

Analog zum arithmetischen Mittel eines Merkmals mit wenig Ausprägungen werden wir aber nicht die obige Formel für die Varianz verwenden, sondern die Vorteile der Häufigkeitstabelle nutzen. Wir können nämlich die Stichprobenvarianz berechnen als

$$s^2 = \frac{1}{n-1} \sum_{j=1}^m (x_j - \bar{x})^2 \cdot n_j,$$

wobei die x_j jetzt die verschiedenen Ausprägungen des Merkmals darstellen.

j	Preis für eine Weinflasche x_j	absolute Häufigkeit n_j	$x_j n_j$	$(x_j - \bar{x})^2$	$(x_j - \bar{x})^2 n_j$
1	4	5	20	5,5696	27,8480
2	5	8	40	1,8496	14,7968
3	7	7	49	0,4096	2,8672
4	10	5	50	13,2496	66,2480
Σ		25	159		111,7600

Zunächst benötigen wir den Mittelwert \bar{x} . Er berechnet sich wie in [Lageparameter](#) als

$$\bar{x} = \frac{1}{n} \sum_{j=1}^m x_j n_j = \frac{159}{25} = 6,36.$$

Wir erhalten nun

$$s^2 = \frac{1}{24} \cdot 111,7600 \approx 4,66.$$

Der Computer kann das leicht ermitteln. Möchten wir jedoch die Varianz händisch ausrechnen, finden wir den „krummen“ Mittelwert als störend. Wir können natürlich auch hier den Verschiebungssatz anwenden. Es gilt nämlich für die benötigte Quadratsumme:

$$Q = \sum_{j=1}^n (x_j - \bar{x})^2 \cdot n_j = \left(\sum_{j=1}^n x_j^2 \cdot n_j \right) - n \cdot \bar{x}^2.$$

Wir berechnen sukzessive in unserer Häufigkeitstabelle die x_j^2 und $x_j^2 n_j$ und erhalten zunächst für Q

$$Q = 1123 - 25 \cdot 6,36^2 = 111,76$$

und für die Varianz

$$s^2 = \frac{111,76}{25 - 1} = 4,66.$$

j	Preis für eine Weinflasche x_j	absolute Häufigkeit n_j	$x_j n_j$	x_j^2	$x_j^2 n_j$
1	4	5	20	16	80
2	5	8	40	25	200
3	7	7	49	49	343
4	10	5	50	100	500
Σ		25	159		1123

4.2.2.

Metrisches Merkmals mit vielen Ausprägungen

<< hoch zum Anfang vom Inhaltsverzeichnis

<< hoch zum Kapitel 4 vom Inhaltsverzeichnis

Varianz als Streuungsparameter

Der Lageparameter allein reicht für die Beschreibung einer Datenmenge nicht aus (analoges Problem wie bei Zufallsverteilungen). Information über die **Streuung** der Beobachtungswerte liefert ein **Streuungsparameter**. Es kommen verschiedene Kennwerte als Streuungsparameter in Betracht, beispielsweise die Varianz, die Spannweite, der Quartilsabstand und der Variationskoeffizient.

Varianz

Am häufigsten wird als Kennwert die Varianz verwendet, da sie wahrscheinlichkeitstheoretisch am besten zu untersuchen ist. Die Varianz sind die mittleren quadratischen Abweichungen der Einzelwerte x_i vom arithmetischen Mittel

$$s^2 = \frac{1}{n - 1} \sum_{i=1}^n (x_i - \bar{x})^2$$

Der Nenner n-1 wirkt vielleicht etwas befremdlich. Allerdings hat die Verwendung von n-1 statt n wahrscheinlichkeitstheoretische Vorzüge, wenn man die Varianz der Verteilung eines Merkmals mit s^2 schätzen möchte. Man nennt dieses Art der Varianz inferentielle Varianz.

Beispiel

Eine Firma möchte einen Kachelofen auf den Markt bringen, der für einen Komplettpreis zu erwerben ist. Für die Kalkulation dieses Preises benötigt die Firma Informationen über die Montagezeit für einen Kachelofen. Bei der Endmontage von 11 Kachelöfen ergaben sich die Zeiten

2,5 3 3 3,3 3,6 3 2,3 3 3,1 3,2 3

Die Varianz der Montagezeiten soll bestimmt werden. Nach der obigen Formel muss zunächst das arithmetische Mittel bestimmt werden:

$$\begin{aligned}\bar{x} &= \frac{1}{11}(2,5 + 3 + 3 + 3,3 + 3,6 + 3 + 2,3 + 3 + 3,1 + 3,2 + 3) \\ &= \frac{33}{11} = 3h\end{aligned}$$

Dann erhalten wir als Varianz

$$\begin{aligned}s^2 &= \frac{1}{10}((2,5 - 3)^2 + (3 - 3)^2 + (3 - 3)^2 + \dots + (3 - 3)^2) \\ &= \frac{1}{10}(0,25 + 0 + 0 + 0,09 + 0,36 + 0 + 0,49 + 0 + 0,01 + 0,04 + 0) \\ &= \frac{1,24}{10} = 0,124h^2\end{aligned}$$

Verzichtet man auf eine Schätzung, kann man auch die deskriptive Varianz

$$s_d^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

für die Beschreibung von statistischen Daten verwenden, was aber hier zur Vermeidung von Verwechslungen unterlassen wird.

Bei der manuellen Berechnung von s^2 ist es oftmals mühsam, erst die einzelnen Differenzen $x_i - \bar{x}$ zu bilden und dann zu quadrieren. Mit Hilfe des **Verschiebungssatzes** kann die laufende Differenzenbildung vermieden werden. Betrachten wir die Summe

$$\sum_{i=1}^n (x_i - \bar{x})^2.$$

Diese Summe lässt sich zerlegen in

$$Q = \sum_{i=1}^n x_i^2 - n \cdot \bar{x}^2.$$

Setzt man den Ausdruck oben ein, erhält man für die Varianz

$$s^2 = \frac{1}{n-1} Q = \frac{1}{n-1} \left(\sum_{i=1}^n x_i^2 - n \cdot \bar{x}^2 \right)$$

Beispiel:

$$\begin{aligned}s^2 &= \frac{1}{10}(2,5^2 + 3^2 + 3^2 + \dots + 3^2 - 11 \cdot 3^2) \\ &= \frac{1}{10} \cdot (100,24 - 99) = 0,124h^2\end{aligned}$$

Da die Varianz ein quadratischer Ausdruck ist, hat sie z.B. auch die Einheit h^2 , wenn die x_i die Einheit h haben. Um die Varianz anschaulicher zu machen, kann man ihre Quadratwurzel, die Standardabweichung s betrachten:

Beispiel

$$s = \sqrt{0,124h^2} \approx 0,35h$$

also ca. 20 Minuten. Man könnte etwas flapsig sagen, dass die Montagezeit eines Ofens im Mittel 3 Stunden +/- 20 Minuten beträgt.

Auch die Varianz reagiert empfindlich auf Ausreißer. Es gibt hier resistente Streuungsparameter, die weiter unten behandelt werden.

4.2.2.1.	Verteilung	<< hoch zum Anfang vom Inhaltsverzeichnis << hoch zum Kapitel 4 vom Inhaltsverzeichnis
4.2.2.1.1.	Klassierung	<< hoch zum Anfang vom Inhaltsverzeichnis << hoch zum Kapitel 4 vom Inhaltsverzeichnis

Metrische Merkmale mit vielen verschiedenen Ausprägungen

Klassierung

Liegen sehr viele verschiedene Beobachtungen eines metrisch skalierten Merkmals vor, ist es wenig sinnvoll, die Ausprägungen zu zählen. Hier müssen die einzelnen Werte für die Häufigkeitstabelle zusammengefasst werden. Das geschieht in sogenannten Klassen.

Beispiel

Es liegen für 32 europäische Länder als Indikator für den Wohlstand die Zahlen der PKWs pro 1000 Einwohner vor:

31	43	65	152	156	247	264	266
280	289	295	332	341	351	357	365
400	421	422	423	438	451	452	456
489	494	514	516	541	557	591	641

Diese Vielzahl unterschiedlicher Werte ist unübersichtlich. Sie werden zu **Klassen** zusammengefasst, und zwar so,

Klasse 1	über 0 - bis 200	31 43 65 152 156
Klasse 2	über 200 bis 300	247 264 266 280 289 295
Klasse 3	über 300 bis 400	332 341 351 357 365 400
Klasse 4	über 400 bis 500	421 422 423 438 451 452 456 489 494
Klasse 5	über 500 bis 700	514 516 541 557 591 641

so dass wir dann die folgende Häufigkeitstabelle erhalten:

j	Zahl der PKW pro 1000	Zahl der Länder absolute Häufigkeit n_j	relative Häufigkeit p_j
1	über 0 - bis 200	5	$5/32 = 0,15625$
2	über 200 bis 300	6	0,1875
3	über 300 bis 400	6	0,1875
4	über 400 bis 500	9	0,28125
5	über 500 bis 700	6	0,1875
Σ		32	1

Struktur von Klassen

Wir wollen anhand des Beispiels die Struktur von Klassen ansehen:

Es werden benachbarte Merkmalsausprägungen x_i zu einer Klasse zusammengefasst. Wir bezeichnen als

- Zahl der Klassen: m ($m=5$)
- Absolute der Beobachtungswerte in der Klasse j ($j = 1, \dots, m$): n_j
- Relative Häufigkeit: $p_j = \frac{n_j}{n}$
- Klassenobergrenze: x_{oj} ; Klassenuntergrenze: x_{uj}
- Klassenbreite: $d_j = x_{oj} - x_{uj}$

- Klassenmitte: $x'_j = \frac{x_{oj} + x_{uj}}{2}$

Bemerkungen

Die Beobachtungen sollen in einer Klasse möglichst gleichmäßig verteilt sein. Idealerweise haben alle Klassen dieselbe Breite, was aber nur bei gleichmäßiger Verteilung der Beobachtung zu empfehlen ist. Auf jeden Fall sollen keine leeren Klassen in der Mitte auftreten.

Für die empfehlenswerte Zahl von Klassen gilt die Faustregel $m \approx \sqrt{n}$. Die Zuordnung der Beobachtung zu einer Klasse muß eindeutig sein, also

nicht	10 - 11	11 - 12	12 - 13	...
sondern	10 - unter 11	11 - unter 12	12 - unter 13	...
oder	über 10 bis 11	über 11 bis 12	über 12 bis 13	...

Manchmal treten **offene Randklassen** auf.

Beispiel:

Größe der landwirtschaftlichen Betriebe in Bayern

Klasse j	Größe des Betriebs(in ha)	...
1	höchstens 5	...
2	über 5 bis 10	...
3	über 10 bis 50	...
4	mehr als 50	...

Behandlung offener Randklassen

Bestimmte Verfahren wie beispielsweise Histogramme etc. verlangen einen Randwert für die oberste und unterste Klasse. Bei offenen Randklassen muß der äußere Randwert „erfunden“ werden.

1. Falls gleiche Klassenbreiten existieren, werden die Randklassen genauso breit gemacht.
2. Man verwendet als äußere Klassengrenze einen plausiblen Wert.

Grafiken

Der Klassiker einer Grafik für klassierte Daten ist das Histogramm, eine Entsprechung des Säulendiagramms. Man trägt auf der Abszisse die Klassen ab und errichtet über den Klassen Rechtecke, deren Fläche die absolute oder relative Häufigkeit beträgt.

Wir wollen nun für die PKW-Indikatordaten ein Histogramm konstruieren. Die Intervallbreiten und die Flächen der einzelnen Rechtecke sind bekannt, uns fehlt jedoch die Höhe einer Säule. Wir werden dieses Problem geometrisch angehen:

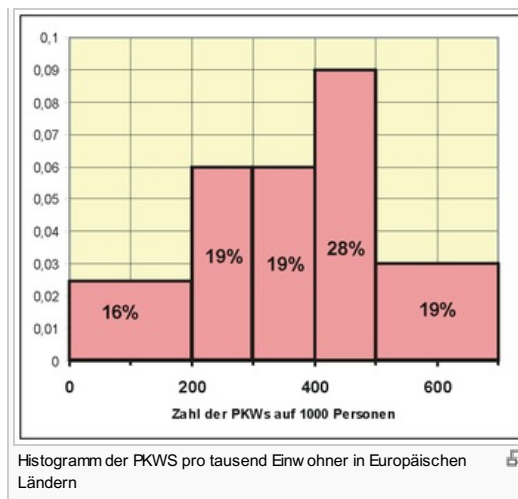
Es gilt Fläche = Höhe * Breite, bzw.

$$n_j = h_j \cdot d_j$$

also

$$h_j = \frac{n_j}{d_j}$$

j	Zahl der PKW pro 1000	Zahl der Länder absolute Häufigkeit n_j	Klassenbreite d_j	Säulenhöhe $h_j = n_j/d_j$
1	über 0 - bis 200	5	200 - 0 = 200	0,025
2	über 200 bis 300	6	100	0,06
3	über 300 bis 400	6	100	0,06
4	über 400 bis 500	9	100	0,09
5	über 500 bis 700	6	200	0,03



Üblicherweise wird beim Histogramm die Ordinate (y-Achse) weggelassen, weil sonst die Höhe der Säule als Häufigkeit gedeutet wird. Tatsächlich ist aber die Fläche der Säule die Häufigkeit. Es geht ja in der Grafik darum, einen optischen Eindruck von der Aufteilung der Daten zu bekommen. In unserem Beispiel wurde die Ordinate beibehalten, damit die Konstruktion des Histogramms deutlich wird. Man kann zur Unterstützung der Information noch die Häufigkeiten in die Säulen eintragen.

4.2.2.1.2

Summenkurve

<< hoch zum Anfang vom Inhaltsverzeichnis

<< hoch zum Kapitel 4 vom Inhaltsverzeichnis

Bei Beobachtungen, die man zweckmäßigerweise **klassiert** zusammenfasst, ist eine Summenfunktion aus der Urliste schwierig zu erstellen und auch unhandlich.

Da hier das Merkmal als stetig angesehen wird, nähert man die grafische Darstellung der Verteilung durch eine Kurve an. Dabei wird folgendermaßen vorgegangen:

Um die absolute Summenfunktion zu erstellen, berechnet man für jede Klasse j die kumulierte Häufigkeit S_j . Dann trägt man die Wertepaare $(x_{oj}; S_j)$, also die Klassenobergrenze und Summenhäufigkeit in ein Diagramm ein und verbindet die Punkte geradlinig. Es ist der erste Punkt $(x_{u1}; 0)$. Ab $(x_{om}; n)$ verläuft die Summenkurve horizontal.

PKW-Beispiel

Dazu fassen wir die benötigten Werte am besten wieder in einer Zahlentabelle zusammen: Wir benötigen die Klassenobergrenzen x_{oj} und die Summenhäufigkeiten S_j . Die Summenhäufigkeiten sind die kumulierten Häufigkeiten

$$S_j = \sum_{k=1}^j n_k$$

etwa $S_1 = 5$, $S_2 = 5 + 6 = 11$, $S_3 = 5 + 6 + 6 = 17 \dots$

Klasse	Merkmalswerte	Absolute Häufigkeit	Klassenobergrenze	Absolute Summenhäufigkeit
j	x	n_j	x_{oj}	S_j

1	0 - bis 200	5	200	5
2	ü. 200 bis 300	6	300	11
3	ü. 300 bis 400	6	400	17
4	ü. 400 bis 500	9	500	26
5	ü. 500 bis 700	6	700	32
Σ		32		

Je gleichmäßiger die einzelnen Beobachtungen über die Klassen verteilt sind, desto besser passt sich die Summenkurve an die Summenfunktion der einzelnen Beobachtungen an.

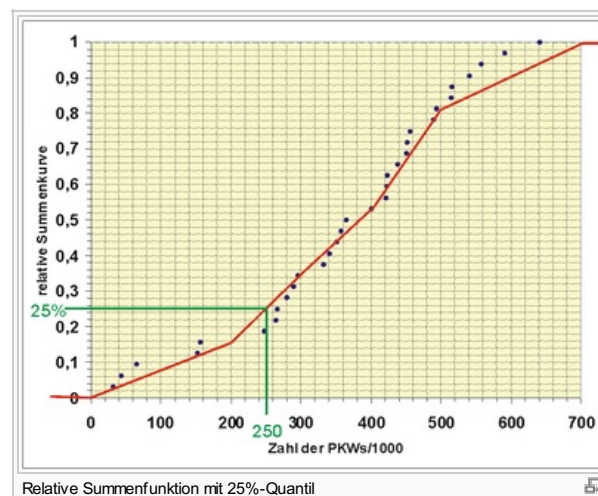
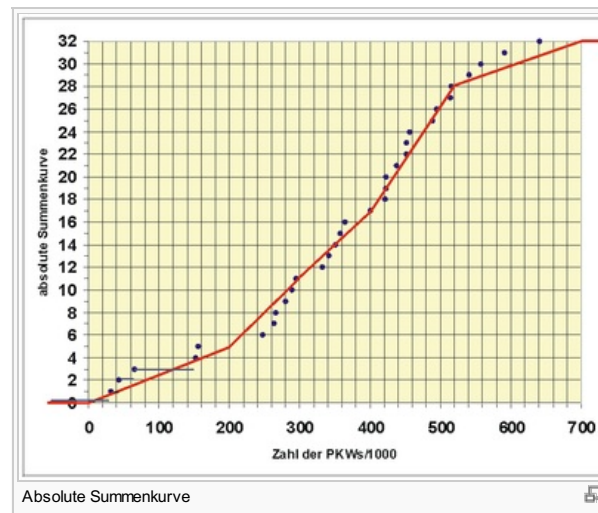
In der Grafik ist die Summenkurve für das PKW-Beispiel angegeben. Zum Vergleich wurde die Summenfunktion der Urliste mit eingetragen, wobei aus Übersichtlichkeitsgründen nur bei den ersten Werten die Horizontale gezeigt wird. Man sieht, dass im Intervall 200 - 300 die Kurve die tatsächlichen Beobachtungen überschätzt, im Intervall 600 - 700 liegt die Kurve unter der tatsächlichen Summenfunktion.

Die Summenfunktion ist eine empirische Beschreibung der Verteilung des Merkmals in der Grundgesamtheit. Wie andere grafisch dargestellte Verteilungen ist auch sie vom optischen Informationsgehalt her eher wenig instruktiv. Man kann aber Verteilungsaussagen grafisch ermitteln, z.B.

Bei der relativen Summenkurve wird statt der absoluten Häufigkeit S_j die relative Summenhäufigkeit

$$S_j^* = \frac{S_j}{n}$$

verwendet. Die Form der Summenkurve bleibt erhalten.



4.2.2.2.

Lageparameter

<< hoch zum Anfang vom Inhaltsverzeichnis

<< hoch zum Kapitel 4 vom Inhaltsverzeichnis

Arithmetisches Mittel

Ist die Urliste gegeben, berechnet sich das arithmetische Mittel aus der bekannten Durchschnittsbildung der Beobachtungswerte. Sind jedoch die Informationen der Urliste nicht mehr verfügbar, kann man das arithmetische Mittel nur noch näherungsweise bestimmen. Man verwendet die Klassenmitte x'_j als Ersatz für die Merkmalsausprägung x_j in der Klasse j und nähert das arithmetische Mittel an als

$$\bar{x} \approx \bar{x}' = \frac{1}{n} \sum_{j=1}^m x'_j \cdot n_j$$

Die Klassenmitte soll das Niveau einer Klasse widerspiegeln. Das ist vor allem der Fall, wenn sich die Einzelwerte der Urliste gleichmäßig in einer Klasse verteilen. Sind die Einzelwerte mehrheitlich an einer Klassengrenze gelegen, gibt x'_j unter Umständen nicht mehr das Niveau korrekt wieder. Die optimale Aufteilung der Klassen sollte schon bei Klassenbildung berücksichtigt werden. Im Sinne einer einfachen Handhabung sollte x'_j eine Zahl sein, mit der man leicht rechnen kann, also z.B. 200 und nicht 199,5.

PKW-Beispiel

Es ergibt sich also als angenähertes arithmetisches Mittel

$$\bar{x}' = \frac{1}{32} \cdot 11750 = 367,1875$$

Klasse	Intervall	Absolute Häufigkeit	Klassenmitte	
j		n_j	x'_j	$x'_j \cdot n_j$
1	0 - bis 200	5	100	500
2	200 bis 300	6	250	1500
3	300 bis 400	6	350	2100
4	400 bis 500	9	450	4050
5	500 bis 700	6	600	3600
Σ		32		11750

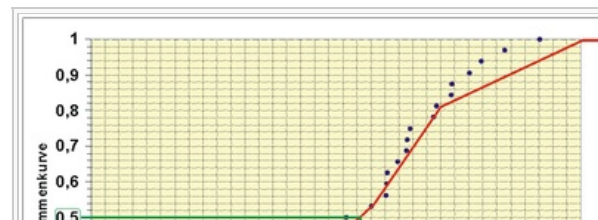
Median

Grafische Ermittlung

Hier bietet sich vor allem die grafische Ermittlung des Medians an:

Man bestimmt aus der absoluten (relativen) Summenkurve grafisch den Wert x , der zu $n/2$ (0,5) gehört.

Im PKW-Beispiel wurde der Median aus der relativen Summenkurve grafisch ermittelt. Der x -Wert, der zu $S^*(X)=0,5$ gehört, beträgt etwa 382. Es hatten also 50% der untersuchten Länder höchstens ca. 382



Fahrzeuge pro 1000 Einwohner.

Ist n klein, könnte man auch vom Ordinatenwert (n+1)/2 bei geradem n ausgehen.

Ermittlung mit der Häufigkeitstabelle

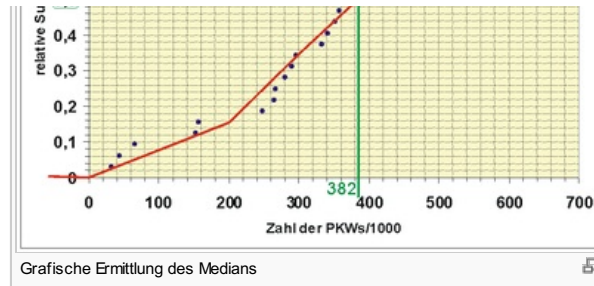
Man kann den Median auch näherungsweise durch lineare Interpolation aus der Häufigkeitstabelle ermitteln. Allerdings genügt im Allgemeinen auch die Klassenmitte der Einfallsklasse als Näherung für den Median, da ohnehin meistens keine Informationen über die Verteilung der Beobachtungen in den Klassen vorliegen.

Im PKW-Beispiel ergäbe die Näherung durch die Klassenmitte $z' = 350$.

Lineare Interpolation würde

$$x_{u3} + \frac{0,5 \cdot (x_{o3} - x_{u3})}{p_3} = 300 + \frac{0,5 \cdot 100}{0,5312} = 394,12$$

ergeben.



4.2.2.3.

Streuungsparameter

[<< hoch zum Anfang vom Inhaltsverzeichnis](#)

[<< hoch zum Kapitel 4 vom Inhaltsverzeichnis](#)

Einleitung

Liegen bei einem klassierten Merkmal keine Informationen über die Urliste mehr vor, können wir die Varianz des Merkmals analog zum arithmetischen Mittel mit den Klassenmitten näherungsweise berechnen. Wir erhalten für die Näherung s^2

$$s^2 \approx s^{2'} = \frac{1}{n-1} \sum_{j=1}^m (x'_j - \bar{x}')^2 \cdot n_j,$$

deren Exaktheit auch wieder von der Verteilung der einzelnen Werte in den Klassen abhängt. Verwenden wir statt der absoluten Häufigkeiten n_j die relativen p_j , berechnet sich die Varianz als

$$s^2 \approx s^{2'} = \frac{n}{n-1} \sum_{j=1}^m (x'_j - \bar{x}')^2 \cdot p_j.$$

Man kann auch im Fall der näherungsweisen Berechnung den Verschiebungssatz anwenden. Wir wollen ihn hier nur für absolute Häufigkeiten angeben. Für die Quadratsumme der zentrierten Klassenmittel gilt

$$\sum_{j=1}^m (x'_j - \bar{x}')^2 \cdot n_j = \sum_{j=1}^m x_j'^2 \cdot n_j - n \cdot \bar{x}'^2,$$

so dass sich für die angenäherte Varianz ergibt

$$s^2 \approx s^{2'} = \frac{1}{n-1} \left(\sum_{j=1}^m x_j'^2 \cdot n_j - n \cdot \bar{x}'^2 \right)$$

PKW-Beispiel

Wie bei der Ermittlung des arithmetischen Mittels verwenden wir auch hier zweckmäßigerweise eine Tabelle. Es war das angenäherte arithmetische Mittel 367, 1875. Es wird zunächst die Varianz mit Hilfe der zentrierten Werte ermittelt:

		Absolute	
--	--	----------	--

Klasse	Intervall	Absolute Häufigkeit	Klassenmitte
j	über ... bis ...	n_j	x_j
1	0 - 200	5	100
2	200 - 300	6	250
3	300 - 400	6	350
4	400 - 500	9	450
5	500 - 700	6	600
Σ	–	32	–

Klasse			
j	$x'_j - \bar{x}_j$	$(x'_j - \bar{x}_j)^2$	$(x'_j - \bar{x}_j)^2 n_j$
1	-267,19	71390,50	356952,48
2	-117,19	13733,50	82400,98
3	-17,19	295,50	1772,98
4	82,81	6857,50	61717,46
5	232,81	54200,50	325202,98
Σ	–	–	828046,88

Wir erhalten für die Varianz

$$s^2 = \frac{1}{32 - 1} \cdot 828046,88 = 26711,19$$

und für die Standardabweichung

$$s = \sqrt{26711,19} = 163,44$$

Mit dem Verschiebungssatz dagegen erhalten wir

Klasse	Intervall	Absolute	Klassen-		
--------	-----------	----------	----------	--	--

	Intervall	Häufigkeit	mitte		
j	über ... bis ...	n _j	x _j '	x _j ' ²	x _j ' ² n _j
1	0 - 200	5	100	10000	50000
2	200 - 300	6	250	62500	375000
3	300 - 400	6	350	122500	735000
4	400 - 500	9	450	202500	1822500
5	500 - 700	6	600	360000	2160000
Σ		32			5142500

Wir erhalten für die Varianz

$$s^2 = \frac{1}{32 - 1} (5142500 - 32 \cdot 367,19^2) = 26711,19$$

5.

Analyse mehrerer Merkmale

[<< hoch zum Anfang vom Inhaltsverzeichnis](#)

[<< hoch zum Kapitel 5 vom Inhaltsverzeichnis](#)

Deskriptive Analyse mehrerer Merkmale

Häufig interessiert man sich für mehrere Merkmale zugleich. Interpretiert man die Beobachtungen wieder als Stichprobe einer unbekannten Grundgesamtheit, könnte man fragen, ob die Variablen unabhängig sind oder, falls nicht, in welcher Beziehung sie zueinander stehen. So kann man beispielsweise etwa vermuten, daß zwischen Werbeausgaben und Umsatz eines Supermarktes ein positiver Zusammenhang besteht.

5.1.

Korrelationsanalyse

[<< hoch zum Anfang vom Inhaltsverzeichnis](#)

[<< hoch zum Kapitel 5 vom Inhaltsverzeichnis](#)

Korrelation zweier Merkmale

Für die Untersuchung der Beziehung zwischen mehreren Variablen muß grundsätzlich wieder nach Skalierung dieser Variablen unterschieden werden. Die Kovarianz bzw. der Korrelationskoeffizient für zwei Zufallsvariablen einer **Grundgesamtheit** sind uns bereits bekannt. Analog dazu gibt es in der deskriptiven Statistik die **(Stichproben)-Kovarianz** bzw. den **(Stichproben)-Korrelationskoeffizienten**.

Korrelationskoeffizient nach Bravais-Pearson

Es seien zwei Merkmale x und y zu beobachten. Bei einer Stichprobe im Umfang von n ergeben sich n viele Wertepaare (x_i; y_i) (i = 1, ..., n).

Beispiel

Es soll untersucht werden, ob das Bevölkerungswachstum eines Landes mit der Fruchtbarkeitsrate (durchschnittliche Zahl der Geburten einer gebärfähigen Frau) zusammenhängt. Es wurden acht Länder zufällig ausgewählt und wir erhalten die Daten

Land	Bevölkerungs- wachstum x	Fruchtbarkeits- rate y

Ägypten	1,8	3
Türkei	1,1	2
Vereinigte Arabische Emirate	1,6	3
Jamaika	0,7	2
Mauritanien	2,9	5
Island	1	1,8
Tadschikistan	2,1	4,1
Gabun	2,4	4,7

Um sich einen Eindruck vom Zusammenhang der Daten zu verschaffen, tragen wir sie in einem Streudiagramm ab.

Man sieht hier schon ohne Analyse, dass offensichtlich mit steigender Fertilität auch das Bevölkerungswachstum zunimmt. Die gestreckte Punktwolke ist fast eine steigende Gerade, also besteht zwischen Fertilität und Bevölkerungswachstum ein annähernd linearer Zusammenhang. Die Merkmale sind offensichtlich stetig. Für metrisch skalierte Merkmale stellt der Korrelationskoeffizient r_{xy} oder kurz r nach Bravais-Pearson ein **Maß für die lineare Abhängigkeit** zweier statistischer Variablen dar:

$$r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \cdot \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}},$$

wobei x_1, x_2, \dots, x_n und y_1, y_2, \dots, y_n die Messwerte der beiden Merkmale

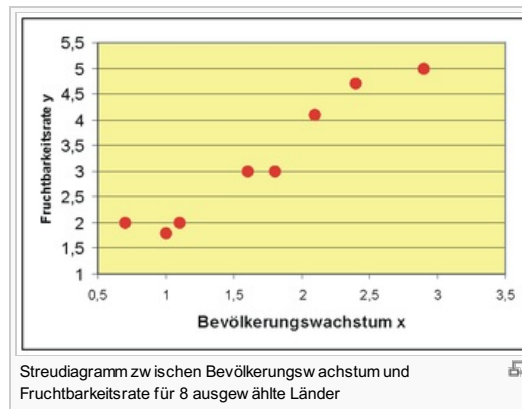
und $\bar{x} = \frac{1}{n} \cdot \sum_{i=1}^n x_i$ das arithmetische Mittel des Merkmals x sind, \bar{y}

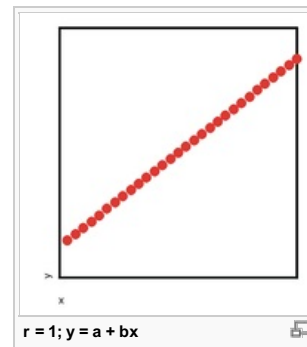
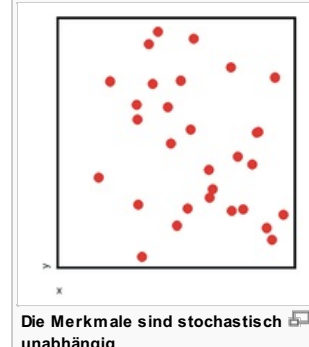
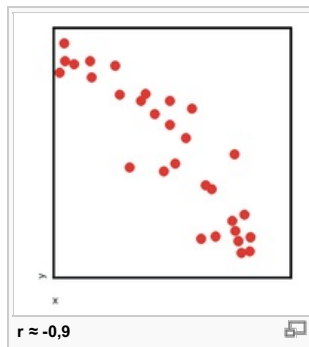
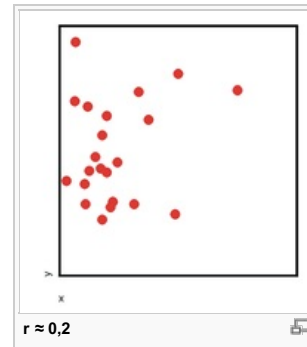
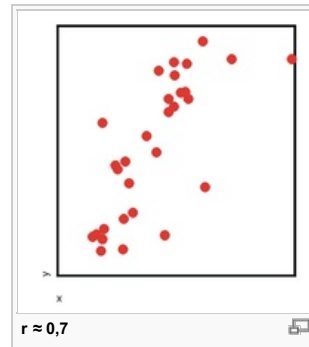
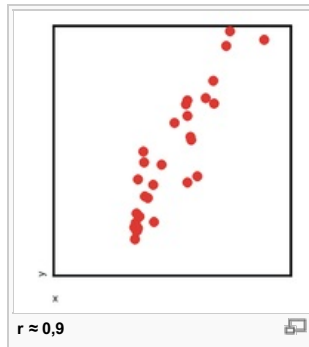
entsprechend.

Analog zu oben kann auch hier wieder der Verschiebungssatz angewendet werden:

$$r = \frac{\sum_{i=1}^n x_i \cdot y_i - n \cdot \bar{x} \cdot \bar{y}}{\sqrt{(\sum_{i=1}^n x_i^2 - n \cdot (\bar{x})^2) \cdot (\sum_{i=1}^n y_i^2 - n \cdot (\bar{y})^2)}}$$

Es gilt: Je näher $|r|$ bei 0 ist, desto schwächer ist der „lineare Zusammenhang“, d.h. die Korrelation. Man sieht an den folgenden Streudiagrammen, dass bei einem Korrelationskoeffizienten von 0,9 das Diagramm stark einer Geraden ähnelt. Je kleiner $|r|$ wird, desto verwaschener wird die Gerade bis hin zur strukturlosen Punktwolke. Ist der Korrelationskoeffizient kleiner als Null, hat die Punktwolke eine fallende Tendenz.





In der Grundgesamtheit ist bei stochastisch unabhängigen Zufallsvariablen die Kovarianz und damit der Korrelationskoeffizient gleich Null. Bei einer Stichprobe stetiger Merkmale wird man aber so gut wie niemals einen Korrelationskoeffizienten erhalten, der genau Null ist. In unserem Beispiel mit den stochastisch unabhängigen Merkmalen wurden 30 Zufallszahlen zweier stochastisch unabhängiger Variablen erzeugt. Der errechnete Stichproben-Korrelationskoeffizient ergab jedoch -0,272. Die Frage ist nun, wie groß muss der errechnete Korrelationskoeffizient mindestens sein, damit man von einer vorhandenen Korrelation ausgehen kann? Hier kann man den Korrelationskoeffizienten statistisch testen, um nachzuprüfen, ob er groß genug ist.

Beispiel mit zentrierten Merkmalswerten

Wir wollen nun den Korrelationskoeffizienten des obigen Beispiels mit der Formel

$$r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \cdot \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$$

ermitteln. Am besten ordnet man die Daten für die Berechnung in einer Tabelle an (siehe unten). Wir benötigen als Erstes den Mittelwert \bar{x} :

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i = \frac{1}{8} \cdot 13,6 = 1,7,$$

entsprechend erhalten wir für y

$$\bar{y} = \frac{1}{8} \cdot 25,6 = 3,2.$$

Wir wollen nun zuerst die Elemente $x_i - \bar{x}$ bestimmen, wir nennen diese zentrierten Werte von x hier x^* :

$$x_1^* = x_1 - \bar{x} = 1,8 - 1,7 = 0,1$$

$$x_2^* = x_2 - \bar{x} = 1,1 - 1,7 = -0,6$$

...

Wir können nun die Formel von oben etwas kürzer schreiben als

$$r = \frac{\sum_{i=1}^n x_i^* \cdot y_i^*}{\sqrt{\sum_{i=1}^n x_i^{*2}} \cdot \sqrt{\sum_{i=1}^n y_i^{*2}}}$$

Setzen wir die entsprechenden Spaltensummen der Tabelle ein, ergibt sich

$$r = \frac{6,47}{\sqrt{3,96 \cdot 11,22}} = 0,9706.$$

Der Korrelationskoeffizient beträgt also 0,9706. x und y sind hochkorreliert: Wenn die Fruchtbarkeitsrate groß ist, wächst die Bevölkerung stark.

	BevW	FrR	$x^* = x - \bar{x}$	$y^* = y - \bar{y}$			
i	x	y	x^*	y^*	x^*y^*	x^{*2}	y^{*2}
1	1,8	3	0,1	-0,2	-0,02	0,01	0,04
2	1,1	2	-0,6	-1,2	0,72	0,36	1,44
3	1,6	3	-0,1	-0,2	0,02	0,01	0,04
4	0,7	2	-1	-1,2	1,2	1	1,44
5	2,9	5	1,2	1,8	2,16	1,44	3,24
6	1	1,8	-0,7	-1,4	0,98	0,49	1,96
7	2,1	4,1	0,4	0,9	0,36	0,16	0,81
8	2,4	4,7	0,7	1,5	1,05	0,49	2,25
Σ	13,6	25,6	0	0	6,47	3,96	11,22

Beispiel mit Verschiebungssatz Wir berechnen Korrelationskoeffizienten mit Hilfe des Verschiebungssatzes:

$$r = \frac{49,99 - 8 \cdot 1,7 \cdot 3,2}{\sqrt{(27,08 - 8 \cdot 1,7^2) \cdot (93,14 - 8 \cdot 3,2^2)}} = 0,9706$$

	BevW	FrR			
i	x	y	xy	x^2	y^2
1	1,8	3	5,4	3,24	9

2	1,1	2	2,2	1,21	4
3	1,6	3	4,8	2,56	9
4	0,7	2	1,4	0,49	4
5	2,9	5	14,5	8,41	25
6	1	1,8	1,8	1	3,24
7	2,1	4,1	8,61	4,41	16,81
8	2,4	4,7	11,28	5,76	22,09
Σ	13,6	25,6	49,99	27,08	93,14

Bemerkungen

- Der Korrelationskoeffizient nach Bravais-Pearson reagiert stark auf Ausreißer in den Beobachtungen. Daher sollten die vorliegenden Daten idealerweise normalverteilten Merkmalen entstammen.
- Aufgrund der Durchschnittsbildung ist er für ordinalskalierte Merkmale nicht zulässig.
- In der praktischen Anwendung werden bei Verwendung des Verschiebungssatzes die Produkte häufig sehr groß. Um Rundungsfehler zu vermeiden, zentriert man hier vor Berechnung des Korrelationskoeffizienten die Datenwerte zu x_i^* und y_i^* wie oben gezeigt.

Rangkorrelationskoeffizient nach Spearman

Für Variablen, die **stark von der Normalverteilung abweichen**, und auch **ordinalskalierte** Variablen, eignet sich der **Rangkorrelationskoeffizient nach Spearman-Pearson**. Hier werden die einzelnen Beobachtungen von x bzw. y der Größe nach geordnet. Jedem Wert wird seine Rangzahl zugewiesen. Es entstehen so n Paare mit Rangzahlen $rg(x_i)$ und $rg(y_i)$. Aus diesen Rängen wird der Korrelationskoeffizient nach Bravais-Pearson errechnet. Man erhält so den Korrelationskoeffizienten nach Spearman-Pearson:

$$r_{SP} = \frac{\sum_i (rg(x_i) - \overline{rg(x)})(rg(y_i) - \overline{rg(y)})}{\sqrt{\sum_i (rg(x_i) - \overline{rg(x)})^2} \sqrt{\sum_i (rg(y_i) - \overline{rg(y)})^2}}.$$

Wenn alle Ränge verschieden sind, kann man die obige Form zu

$$r_{SP} = 1 - \frac{6 \sum_i d_i^2}{n \cdot (n^2 - 1)},$$

umformen mit $d_i = rg(x_i) - rg(y_i)$.

Liegen mehrere gleiche Merkmalswerte vor, handelt es sich um **Bindungen**. Die untere der beiden Formeln ist eigentlich nur korrekt anwendbar, wenn keine Bindungen vorliegen. Meistens kann man jedoch zur Vereinfachung die Formel näherungsweise verwenden. Zur konkreten Berechnung von Bindungen soll das folgende Beispiel verwendet werden.

Beispiel: Evaluation einer Vorlesung

Es wurde eine Statistikvorlesung evaluiert. Die gesamten Daten sind unter [Evaluation](#) verfügbar. Es wurden hier 10 Studierende zufällig ausgewählt. Wir interessieren uns für die Frage, ob möglicherweise die Zufriedenheit der Leute mit der Vorlesung davon abhängt, ob die Vorlesung verständlich war. Es ergaben sich die Daten

Stoff verständlich	Note für Vorlesung
-----------------------	-----------------------

x	y
2	1
4	4
2	2
3	3
4	3
3	2
3	2
4	3
3	3
3	3

Es werden nun die Ränge ermittelt. Da mehrere Merkmalswerte gleich sind, liegen Bindungen vor, d.h. gleiche Werte bekommen gleiche Rangzahlen. Es gibt verschiedene Methoden, gleiche Rangzahlen zuzuweisen. Meistens werden mittlere Rangzahlen verwendet. Wir wollen für x die Rangzahlen ermitteln. Dazu ordnen wir die x-Werte der Größe nach und numerieren sie durch:

x aufsteigend geordnet	Laufende Nummer	mittlerer Rang	Rangzahl
2	1	$\frac{1+2}{2}$	1,5
2	2		1,5
3	3	$\frac{3+4+5+6+7}{5}$	5
3	4		5
3	5		5
3	6		5
3	7		5
4	8	$\frac{8+9+10}{3}$	9

4	9		9
4	10		9

Für die Ränge von y verfahren wir entsprechend, wie die unten folgende Tabelle zeigt. Nun können wir den Korrelationskoeffizienten nach Spearman-Pearson berechnen:

$$r_{SP} = \frac{\sum_i (rg(x_i) - \overline{rg(x)})(rg(y_i) - \overline{rg(y)})}{\sqrt{\sum_i (rg(x_i) - \overline{rg(x)})^2} \sqrt{\sum_i (rg(y_i) - \overline{rg(y)})^2}}$$

$$= \frac{54,5}{\sqrt{70} \cdot \sqrt{70,5}} = 0,7758,$$

wobei sich für $\overline{rg(x)} = \frac{1}{10} \cdot 55 = 5,5$ ergibt, für $\overline{rg(y)}$ ebenfalls. Es scheint zwischen dem Verstehen des Statistikstoffs und der

Gesamtzufriedenheit ein deutlich positiver Zusammenhang zu bestehen: Je besser der Stoff verstanden wurde, desto besser fiel tendenziell auch die Note aus.

x	y	rg(x)	rg(y)	$rg(x)^* = \frac{rg(x) - \overline{rg(x)}}{\overline{rg(x)} - \overline{rg(x)}}$
2	1	1,5	1	-4
4	4	9	10	3,5
2	2	1,5	3	-4
3	3	5	7	-0,5
4	3	9	7	3,5
3	2	5	3	-0,5
3	2	5	3	-0,5
4	3	9	7	3,5
3	3	5	7	-0,5
3	3	5	7	-0,5
		55	55	0

x	y	$rg(y)^* = \frac{rg(y) - \overline{rg(y)}}{\overline{rg(y)} - \overline{rg(y)}}$	$rg(x)^* \cdot rg(y)^*$	$rg(x)^{*2}$	$rg(y)^{*2}$
2	1	-4,5	18	16	20,25

4	4	4,5	15,75	12,25	20,25
2	2	-2,5	10	16	6,25
3	3	1,5	-0,75	0,25	2,25
4	3	1,5	5,25	12,25	2,25
3	2	-2,5	1,25	0,25	6,25
3	2	-2,5	1,25	0,25	6,25
4	3	1,5	5,25	12,25	2,25
3	3	1,5	-0,75	0,25	2,25
3	3	1,5	-0,75	0,25	2,25
		0	54,5	70	70,5

Wir werden nun den Korrelationskoeffizienten zum Vergleich mit der vereinfachten Formel ermitteln:

$$r_{SP} = 1 - \frac{6 \sum_i d_i^2}{(n \cdot (n^2 - 1))} = 1 - \frac{6 \cdot 31,5}{10 \cdot (100 - 1)} = 0,8091.$$

Dieser Wert weicht etwas vom vorhergehenden ab.

x	y	rg(x)	rg(y)	d _i = rg(x)-rg(y)	d _i ²
2	1	1,5	1	0,5	0,25
4	4	9	10	-1	1
2	2	1,5	3	-1,5	2,25
3	3	5	7	-2	4
4	3	9	7	2	4
3	2	5	3	2	4
3	2	5	3	2	4
4	3	9	7	2	4
3	3	5	7	-2	4

3	3	5	7	-2	4
					31,5

Bemerkungen

- Wie beim Korrelationskoeffizienten nach Bravais-Pearson kann auch hier der Verschiebungssatz verwendet werden.
- Wird für die Berechnung des Korrelationskoeffizienten der Computer eingesetzt, sollte die vereinfachte Formel nicht verwendet werden, denn sie soll lediglich bei der Berechnung von Hand die Arbeit erleichtern - es sei denn, alle Rangzahlen sind verschieden.

5.2.

Regressionsanalyse

<< hoch zum Anfang vom Inhaltsverzeichnis

<< hoch zum Kapitel 5 vom Inhaltsverzeichnis

Einfaches lineares Regressionsmodell Einführung mit Beispiel einer Preis-Absatz-Funktion

Eine renommierte Sektkellerei möchte einen hochwertigen Rieslingsekt auf den Markt bringen. Für die Festlegung des Abgabepreises soll zunächst eine Preis-Absatz-Funktion ermittelt werden. Dazu wurde in $n = 6$ Geschäften ein Testverkauf durchgeführt. Man erhielt sechs Wertepaare mit dem Ladenpreis x (in Euro) einer Flasche und die verkaufte Menge y an Flaschen:

Laden	i	1	2	3	4	5	6
Preis einer Flasche	x_i	20	16	15	16	13	10
verkaufte Menge	y_i	0	3	7	4	6	10

Modell

Man geht von folgendem statistischen Modell aus:

Man betrachtet zwei Variable, die vermutlich ungefähr in einem linearen Zusammenhang

$$y \approx \alpha + \beta x$$

stehen. Dabei sind x als unabhängige und y als abhängige Variable definiert. Man nennt auch x erklärende oder exogene Variable und y Zielvariable oder endogene Variable. Es existieren von x und y je n Beobachtungen x_i und y_i ($i = 1, \dots, n$). Der funktionale Zusammenhang $y = f(x)$ zwischen x und y kann nicht exakt festgestellt werden, da $\alpha + \beta x$ von einer Störgröße u überlagert wird, die nichterfassbare Einflüsse (menschliches Verhalten, Messungenauigkeiten usw.) mit einschließt. Es ergibt sich also das Modell

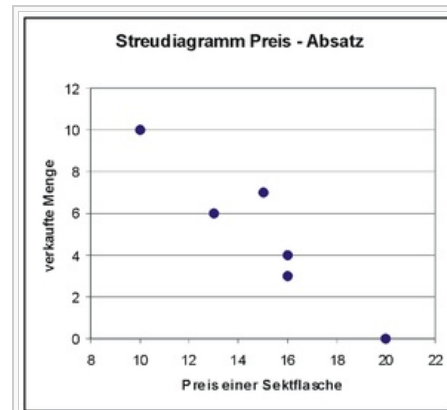
$$y = \alpha + \beta x + u$$

mit den einzelnen Beobachtungen

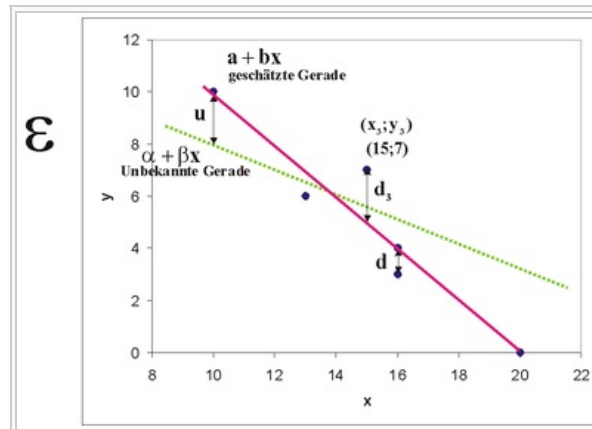
$$y_i = \alpha + \beta x_i + u_i.$$

Da α und βx nicht bekannt sind, kann y auch nicht in die Komponenten $\alpha + \beta x$ und u zerlegt werden.

Es soll eine mathematische Schätzung für die Parameter α und β



Streudiagramm von Preis und abgesetzter Menge an Sektflaschen



Wahre unbekannte und geschätzte Regressionsgerade

durch zwei Konstanten a und b gefunden werden, und zwar so, daß sich ergibt

$$y_i = a + bx_i + d_i.$$

wobei d_i das Residuum bezeichnet, die Abweichung des beobachteten y - Wertes vom geschätzten. Es gibt verschiedene Möglichkeiten, die Regressionsgerade zu schätzen. Man könnte eine Gerade so durch den Punkteschwarm legen, dass die Quadratsumme der Residuen, also der senkrechten Abweichungen d_i der Punkte von dieser **Ausgleichsgeraden** minimiert wird.

Beispiel zum Modell der Grundgesamtheit

In einem breit angelegten Versuch wird ein Flüssigdünger an in Nährlösung gezogenen Peperonis untersucht. Es wird wöchentlich jeder Pflanze eine bestimmte Menge Dünger verabreicht. Nach zwei Monaten wird der Gesamtertrag einer Pflanze gewogen. Die Abhängigkeit des Ertrags y (kg) von der Düngermenge x (ml) lässt sich beschreiben als

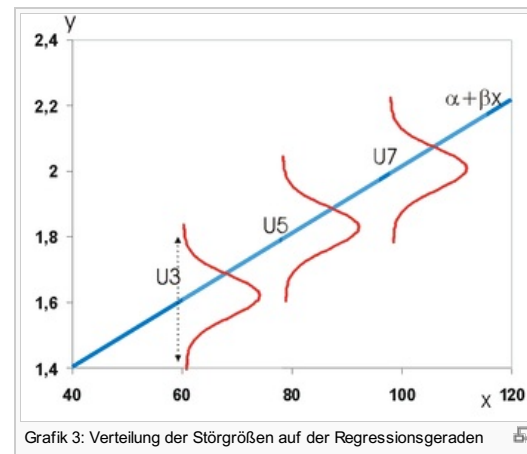
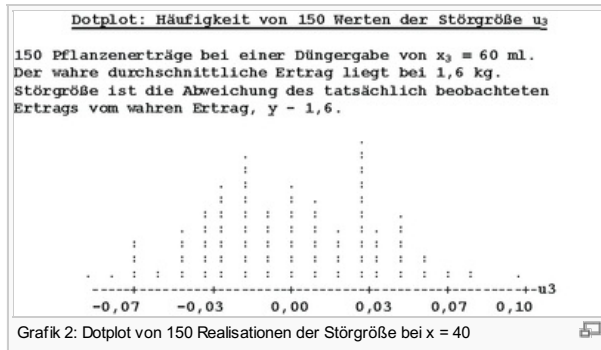
$$y = \alpha + \beta x + u = 1 + 0,01x + u,$$

wobei natürlich kein Mensch diese Gerade kennt, wir können nur einzelne Versuche machen. In der Tabelle 1 sind für die Düngergaben 40, 50, ..., 120 ml für jeweils 8 Pflanzen die resultierenden Erträge aufgeführt. Man sieht, dass die Erträge um $\alpha + \beta \cdot x$ schwanken, was natürlich an der Störgröße $u = y - (\alpha + \beta \cdot x)$ liegt.

Betrachten wir die Störgröße bei einer Düngermenge von $x_3 = 60$ ml. Es wurden hier die Erträge von 150 Peperoni-Pflanzen erfasst. Wenn man die Realisationen der Störgröße u_3 in einem Dotplot abträgt (Grafik 2), erkennt man, dass die Werte normalverteilt sein könnten. Zu x_3 gehört also eine eigene Wahrscheinlichkeitsverteilung der Störgröße, ebenso zu x_1, x_2 usw. In der Grafik 3 sind diese verschiedenen Verteilungen der u exemplarisch angedeutet.

Wöchentliche Düngergabe pro Pflanze (ml)									
	x1	x2	x3	x4	x5	x6	x7	x8	x9
	40	50	60	70	80	90	100	110	120
Ertrag (kg)	1,41	1,47	1,45	1,70	1,57	1,77	2,01	2,32	2,26
	1,44	1,56	1,64	1,58	1,56	2,11	1,79	2,02	2,30
	1,24	1,28	1,62	1,71	1,79	1,92	2,09	2,13	2,21
	1,22	1,35	1,40	1,46	1,78	1,91	2,08	2,16	2,33
	1,26	1,30	1,80	1,80	1,74	2,14	1,79	2,25	2,39
	1,37	1,33	1,57	1,83	1,59	1,90	1,81	1,91	2,18
	1,18	1,36	1,59	1,70	2,03	1,75	2,08	2,23	2,33
	1,56	1,49	1,60	1,75	1,74	1,87	2,13	2,07	1,99
$\alpha + \beta x$	1,4	1,5	1,6	1,7	1,8	1,9	2	2,1	2,2
Störgröße $u = y - (\alpha + \beta x)$	u1	u2	u3	u4	u5	u6	u7	u8	u9
	0,01	-0,03	-0,15	0,00	-0,23	-0,13	0,01	0,22	0,06
	0,04	0,06	0,04	-0,12	-0,24	0,21	-0,21	-0,08	0,10
	-0,16	-0,22	0,02	0,01	-0,01	0,02	0,09	0,03	0,01
	-0,18	-0,15	-0,20	-0,24	-0,02	0,01	0,08	0,06	0,13
	-0,14	-0,20	0,20	0,10	-0,06	0,24	-0,21	0,15	0,19
	-0,03	-0,17	-0,03	0,13	-0,21	0,00	-0,19	-0,19	-0,02
	-0,22	-0,14	-0,01	0,00	0,23	-0,15	0,08	0,13	0,13
	0,16	-0,01	0,00	0,05	-0,06	-0,03	0,13	-0,03	-0,21

Tabelle 1: Daten



Axiome des linearen Regressionsmodells

Damit dieses Verfahren sinnvolle Ergebnisse liefert, wurden für das Lineare Regressionsmodell bestimmte verteilungstheoretische Annahmen

getroffen. Wir gehen aus von der Beziehung

$$y_i = \alpha + \beta x_i + u_i.$$

und definieren die Störgröße u_i als Zufallsvariable. Die Annahmen des linearen Regressionsmodell sind

1. Alle u_i haben den Erwartungswert Null: $Eu_i = 0$, ($i = 1, \dots, n$).
2. Alle u_i haben die gleiche Varianz (Homoskedastizität): $\text{var}u_i = \text{var}u_j$ ($i, j = 1, \dots, n, i \neq j$).
3. Die u_i sind sämtlich stochastisch unabhängig voneinander.

Der Sinn dieser Axiome wird weiter unten erläutert.

Minimierung

Die herkömmliche Methode, die sich auf der Basis der Axiome ergibt, ist die Minimum-Quadrat-Methode oder [Methode der kleinsten Quadrate](#).

Man minimiert also die summierten Quadrate der Residuen,

$$RSS = \sum_{i=1}^n d_i^2 = \sum_{i=1}^n (y_i - (a + bx_i))^2 \rightarrow \min!$$

bezüglich a und b .

Wir multiplizieren die Klammer aus:

$$\begin{aligned} S &= \sum_{i=1}^n (y_i - a - bx_i) \cdot (y_i - a - bx_i) \\ &= \sum_{i=1}^n (y_i^2 - y_i a - y_i bx_i - ay_i + a^2 + abx_i - y_i bx_i + abx_i + b^2 x_i^2) \\ &= \sum_{i=1}^n (y_i^2 - 2y_i a - 2y_i bx_i + a^2 + 2abx_i + b^2 x_i^2) \\ &= \sum_{i=1}^n y_i^2 - 2a \sum_{i=1}^n y_i - 2b \sum_{i=1}^n y_i x_i + na^2 + 2ab \sum_{i=1}^n x_i + b^2 \sum_{i=1}^n x_i^2. \end{aligned}$$

Wir minimieren durch Ableiten

$$\begin{aligned} \frac{\partial S}{\partial a} &= -2 \sum_{i=1}^n y_i + 2na + 2b \sum_{i=1}^n x_i, \\ \frac{\partial S}{\partial b} &= -2 \sum_{i=1}^n x_i y_i + 2a \sum_{i=1}^n x_i + 2b \sum_{i=1}^n x_i^2, \end{aligned}$$

und Nullsetzen, was ein wenig optisch geschönt die Normalgleichungen

$$\begin{aligned} na + b \sum_{i=1}^n x_i &= \sum_{i=1}^n y_i, \\ a \sum_{i=1}^n x_i + b \sum_{i=1}^n x_i^2 &= \sum_{i=1}^n x_i y_i \end{aligned}$$

ergibt.

Wir erhalten die gesuchten Regressionskoeffizienten als die Lösungen

$$b = \frac{\sum_{i=1}^n x_i y_i - n \bar{x} \bar{y}}{\sum_{i=1}^n x_i^2 - n \bar{x}^2}$$

und

$$a = \bar{y} - b \bar{x},$$

wobei $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$ der Mittelwert, also der Durchschnitt der x -Daten ist, y entsprechend. Wegen des Verschiebungssatzes kann man b

auch darstellen als

$$b = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

oder, nach Erweiterung des Bruchs durch $1/(n-1)$,

$$b = \frac{s_{xy}}{s_x^2}$$

mit s_{xy} als **Kovarianz** zwischen den x_i und y_i und s_x^2 als **Varianz** der x_i . Man nennt diese Schätzungen auch **Kleinste-Quadrate-Schätzer**, **KQ**- oder **OLS-Schätzer**.

Wir wollen nun für das obige Sektbeispiel die Regressionskoeffizienten bestimmen:

Preis einer Flasche	verkaufte Menge	$x_i - \bar{x}$	$y_i - \bar{y}$				
x_i	y_i	x^*	y^*	x^*y^*	x^*x^*	y^*y^*	\hat{y}
20	0	5	-5	-25	25	25	0,09
16	3	1	-2	-2	1	4	4,02
15	7	0	2	0	0	4	5,00
16	4	1	-1	-1	1	1	4,02
13	6	-2	1	-2	4	1	6,96
10	10	-5	5	-25	25	25	9,91
90	30	0	0	-55	56	60	30,00

Wir berechnen in dem Beispiel zunächst die arithmetischen Mittel als die Koordinaten des Schwerpunktes der n Messwerte bzw. der Punktwolke.

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i = \frac{1}{6} 90 = 15, \text{ entsprechend } \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i = \frac{1}{6} 30 = 5,$$

und dann die Regressionskoeffizienten

$$b = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{-55}{56} = -0,98 \text{ als die Steigung der Regressionsgeraden}$$

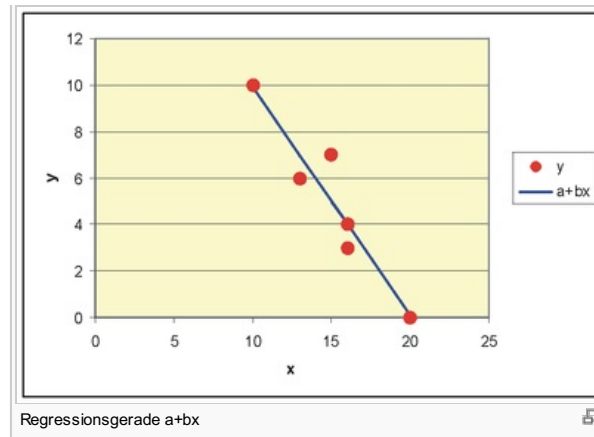
und

und

$$a = \bar{y} - b\bar{x} = 5 + 0,98 \cdot 15 = 19,7$$

Die geschätzte Regressionsgerade lautet $\hat{y} = 19,73 - 0,98x$, so dass man vermuten kann, dass bei jedem Euro mehr der Absatz im Durchschnitt um ca. 1 Flasche sinkt.

Für die gegebenen x-Werte erhalten wir als Schätzungen \hat{y}



$$\hat{y}_1 = a + bx_1 = 19,73 - 0,98 \cdot 20 = 0,09$$

$$\hat{y}_2 = a + bx_2 = 19,73 - 0,98 \cdot 16 = 4,02$$

\vdots

$$\hat{y}_6 = a + bx_6 = 19,73 - 0,98 \cdot 10 = 9,91$$

Für die beobachteten Absatzwerte y bleibt das Residuum r_i übrig:

$$\begin{aligned} y_1 &= a + bx_1 + d_1 = \hat{y}_1 + d_1 \\ \Rightarrow d_1 &= y_1 - \hat{y}_1 = 0 - 0,09 = -0,09 \end{aligned}$$

$$\begin{aligned} y_2 &= a + bx_2 + d_2 = \hat{y}_2 + d_2 \\ \Rightarrow d_2 &= y_2 - \hat{y}_2 = 3 - 4,02 = -1,02 \end{aligned}$$

\vdots

$$\begin{aligned} y_6 &= a + bx_6 + d_6 = \hat{y}_6 + d_6 \\ \Rightarrow d_6 &= y_6 - \hat{y}_6 = 10 - 9,91 = 0,09 \end{aligned}$$

Schätzung der Varianzen

Die Stichprobenvarianz der Residuen berechnet sich als:

$$s^2 = \frac{1}{n-2} \sum_i d_i^2$$

Man schätzt damit die Varianz der Störgröße u (eigentlich U !).

Gesetzmäßigkeiten

Bezüglich der Zielvariablen und der Residuen gilt:

- $\sum_i d_i = 0$

und damit $\bar{d} = 0$.

Die Residuen sind im Mittel Null, sie enthalten also keine Information mehr.

$$\sum_i x_i d_i = 0$$

Die unabhängige Variable x und die Residuen sind orthogonal. Geometrisch bedeutet das, dass sie senkrecht aufeinander stehen. Sie sind daher unkorreliert. Die Residuen enthalten also keinerlei Information mehr, die in x steckt. Die Informationen aus x sind alle in $a + bx$. Nur, was von x nicht mehr erklärt werden kann, bleibt als Rest in d .

$$\bar{\hat{y}} = \frac{1}{n} \sum_i (y_i - d_i) = \frac{1}{n} \sum_i y_i - \frac{1}{n} \sum_i d_i = \bar{y}.$$

Das arithmetische Mittel der \hat{y} ist gleich dem Mittel von y .

Vorteile der Minimum-Quadrat-Methode:

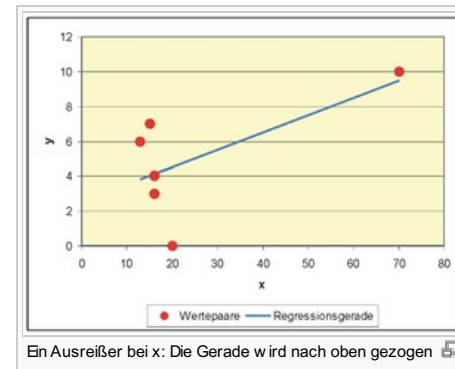
- Positive und negative Abweichungen heben sich bei Summenbildung nicht auf.
- Große Residuen werden im Verhältnis stärker gewichtet als kleine.
- Der Durchschnitt der Residuen ist Null.
- Die Regressionskoeffizienten können mit einer Formel berechnet werden.

Nachteil der Minimum-Quadrat-Methode:

Nicht resistent gegenüber Ausreißern in den Daten.

Sekt-Beispiel:

Wegen eines Erhebungsfehlers wurde für x_6 statt 10 der Wert 70 eingetippt. Die neue Regressionsgerade ergibt sich als $\hat{y} = 2,51 + 0,10x$. Dieser Ausreißer beeinträchtigt das Ergebnis so sehr, dass sogar das Vorzeichen der Steigung umgedreht wird. Eigentlich sollte die Regressionsgerade durch die Punktwolke auf der linken Seite der Grafik führen und fallend sein. Der Ausreißer hebt die Gerade regelrecht aus: Man spricht von einem High-Leverage-Value, also einem Wert mit großer Hebelkraft. Wir erkennen sofort, dass dieser Ausreißer die Analyse völlig wertlos gemacht hat. In dem speziellen Sachzusammenhang könnte man sogar einen fatalen Fehlschluss machen: Bei Luxusgütern sind steigende Preis-Absatz-Funktionen denkbar, weil ein hoher Preis statusfördernd ist. Man könnte also fälschlicherweise annehmen, dass dieser Zusammenhang auch hier gilt. Man würde also einen sehr hohen Preis festlegen und am Markt scheitern.



Bestimmtheitsmaß

Ein Kriterium für die Beurteilung der Güte einer Regressionsschätzung ist das Bestimmtheitsmaß. Die Begründung für dieses Maß leitet sich aus der sog. Streuungszerlegung her. Die Gesamtvarianz von y lässt sich, ausgehend von der Beziehung

$$y_i = \hat{y}_i + d_i$$

zerlegen in die durch $a + bx$ erklärte Varianz von y und die nicht erklärte Varianz:

$$s_y^2 = \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 + \frac{1}{n} \sum_{i=1}^n (d_i - \bar{d})^2.$$

Einige Umformungen ergeben das Bestimmtheitsmaß

$$r^2 = \frac{\frac{1}{n} \sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{\frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2}$$

als Anteil der erklärten Streuung an der Gesamtstreuung von y . Es ist

$$r^2 = \frac{(\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}))^2}{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (y_i - \bar{y})^2},$$

wobei ersichtlich ist, daß r^2 das Quadrat des Korrelationskoeffizienten von x und y darstellt. Mit dem Verschiebungssatz erhalten wir

$$r^2 = \frac{(\sum_{i=1}^n x_i y_i - n \cdot \bar{x} \cdot \bar{y})^2}{(\sum_{i=1}^n x_i^2 - n \cdot \bar{x}^2)(\sum_{i=1}^n y_i^2 - n \cdot \bar{y}^2)}.$$

Es gilt:

$$0 \leq r^2 \leq 1$$

Je näher r^2 bei 1 ist, desto größer ist der Anteil der erklärten Streuung, desto besser wird y durch x erklärt. $r^2 = 0$ bedeutet, dass x und y unkorreliert sind, und $r^2 = 1$, dass x und y eine Gerade bilden.

Die Berechnung der Varianz der Residuen von Hand mit der Formel

$$s^2 = \frac{1}{n-2} \sum_i d_i^2$$

ist aufwendig, weil zuerst die Residuen ermittelt werden müssen. Eine vereinfachte Form leitet sich aus der Beziehung

$$r^2 = 1 - \frac{\sum_{i=1}^n d_i^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

her. Es ist dann nämlich

$$s^2 = \frac{1}{n-2} (1 - r^2) \cdot \sum_{i=1}^n (y_i - \bar{y})^2$$

Sekt-Beispiel

Da hier die arithmetischen Durchschnitte glatte Werte sind, wollen wir das Bestimmtheitsmaß mit der Formel

$$r^2 = \frac{(\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}))^2}{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (y_i - \bar{y})^2},$$

ermitteln. Die Quadratsummen wurden oben in der Tabelle bereits ausgerechnet. Wir erhalten

$$r^2 = \frac{-55^2}{56 \cdot 60} = 0,9003.$$

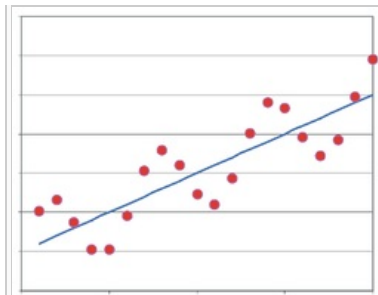
Man könnte also sagen, dass etwa 90% der Information in y von x stammen, die restlichen 10% haben andere Ursachen.

Anforderungen an das Regressionsmodell

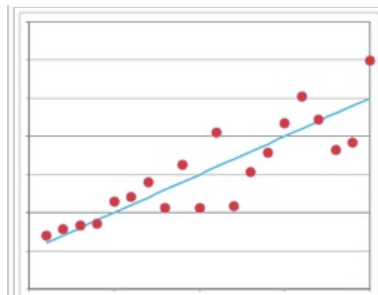
Das Regressionsmodell kann nur optimale Ergebnisse liefern, wenn bestimmte Anforderungen erfüllt sind. Diese Anforderungen lassen sich aus dem Axiomensystem des klassischen linearen Regressionsmodells herleiten:

Die Residuen sollen nur rein zufällig streuen und keinerlei Systematik mehr enthalten, d.h. die Zielvariable y soll durch x vollständig erklärt werden. Systematik in den Residuen deutet daraufhin, daß das Modell möglicherweise falsch bestimmt wurde, d.h. es liegt ein Spezifikationsfehler vor.

Als bestes Mittel zur Überprüfung dieser Modellvoraussetzungen wird das (x;y)-Streudiagramm angesehen, das schnell einen optischen Eindruck von der Verteilung der Störgröße vermittelt.



Korrelierte Residuen:
In den Residuen ist noch ein Schw ingungskomponente, die man ev. mit dem Ansatz $y = a + b_1x + b_2\sin(x)$ einbinden könnte.



Verschiedene Varianz der Residuen:
Die linken Residuen schw anken schw ächer als die rechten. Vermutlich sind zw ei verschiedene Populationen gemischt w orden.

Prognose

Ein Ziel der Regressionsanalyse ist die Prognose \hat{y}_0 , d.h. man fragt danach, welchen Wert y annimmt, wenn ein bestimmtes x_0 vorgegeben ist:

$$\hat{y}_0 = a + bx_0$$

Sekt-Beispiel: Wieviel Flaschen Sekt werden im Durchschnitt verkauft, wenn der Preis auf $x_0 = 11$ Euro festgelegt wird? Es ergibt sich der Prognosewert

$$\hat{y}_0 = 19,7321 - 0,9821 \cdot 11 = 8,93$$

Das heißt jetzt aber nicht, dass in jedem Laden genau 8,93 Flaschen verkauft werden, was auch schwierig sein dürfte, sondern dass in einem Laden durchschnittlich 8,93 Flaschen abgesetzt werden.

Je weiter x_0 vom „Zentrum“ \bar{x} der Daten entfernt ist, desto unverlässlicher werden die Prognosen - ihre Varianz wird immer größer. Deshalb sollte man sich bei einer Prognose nicht zu weit von den Daten entfernen.

Multiple Regression

Beispiel mit demografischen Daten ausgewählter Länder:

Row i	Country	popgrow	birth	mort	fertil	explife	infmort
		y	x_1	x_2	x_3	x_4	x_5
1	Austria	0,14	8,90	9,56	1,35	78,87	4,68
2	Burkina Faso	2,57	44,46	18,79	6,28	44,20	98,67
3	Burma	0,47	18,64	12,16	2,08	56,01	68,78
4	Chile	1,01	15,77	5,71	2,06	76,38	9,05
5	Costa Rica	1,52	18,99	4,32	2,33	76,63	10,26
6	Iraq	2,74	33,09	5,66	4,40	68,26	52,71
7	Norway	0,41	11,89	9,51	1,78	79,25	3,73
8	Portugal	0,41	10,90	10,37	1,46	77,35	5,13
9	Singapore	1,71	9,63	4,05	1,04	81,53	2,28
10	United Kingdom	0,29	10,88	10,19	1,66	78,27	5,22

Erklärung der Variablen:

birth	Geburtenrate (Zahl der Geburten pro 1000 Einwohner)
explife	Lebenserwartung (Mittlere Zahl der Lebensjahre von Geburt an)
fertil	Fertilität (Durchschn. Kinderzahl pro gebärfähiger Frau)
infmort	Kindersterblichkeit (Zahl der verstorbenen Kinder unter einem Jahr pro 1000 Lebendgeburten)
mort	Sterblichkeit (Zahl der Verstorbenen pro 1000 Einwohner)
popgrow	Wachstumsrate der Bevölkerung (prozentuale Änderung der Einwohnerzahl gegenüber dem Vorjahr)

Es wurden die demografischen Daten für n=10 zufällig ausgewählte Länder erhoben (Quelle: [Worldfact-Book der CIA](#) 📄)

Es soll nun das Bevölkerungswachstum *popgrow* erklärt werden. Es wird zunächst als erklärende Variable die Geburtenrate *birth* versucht:

$$popgrow = a + b \cdot birth \text{ bzw. } y = a + bx .$$

Wir erhalten die Regressionsgerade

$$popgrow = -0,104 + 0,0672 \cdot birth$$

mit einem Bestimmtheitsmaß von 66,4%. Die Information in *popgrow* wird also zu 66% durch *birth* erklärt, die restlichen 34% entstammen anderen Einflussgrößen. Wir machen einen zweiten Versuch und verwenden die Sterblichkeit als erklärende Variable:

$$popgrow = a + b \cdot mort + d .$$

Hier ergibt sich als Regressionsgerade

$$popgrow = 1,16 - 0,0032 \cdot mort + d$$

mit einem Bestimmtheitsmaß von ca. 0%. Dieses Ergebnis ist enttäuschend und auch das vorherige war nicht gerade berauschend. Jetzt versuchen wir mal was Innovatives: Wir machen einen Regressionsansatz mit zwei unabhängigen Variablen

$$popgrow = b_0 + b_1 \cdot birth + b_2 \cdot mort + d \text{ bzw. } y = b_0 + b_1 \cdot x_1 + b_2 \cdot x_2 + d .$$

Gesucht ist also die geschätzte **Regressionebene**

$$\hat{y} = b_0 + b_1 x_1 + b_2 x_2$$

Wir erhalten das Gleichungssystem

$$y_1 = b_0 + b_1 x_{11} + b_2 x_{12} + d_1 ,$$

$$y_2 = b_0 + b_1 x_{21} + b_2 x_{22} + d_2 ,$$

$$y_3 = b_0 + b_1 x_{31} + b_2 x_{32} + d_3 ,$$

$$\vdots$$
$$\hat{y}_{10} = b_0 + b_1 x_{10,1} + b_2 x_{10,2} + d_{10} .$$

Wir wollen nun die einzelnen Daten zu Matrizen zusammenfassen. Wir erhalten die (10x3)-Datenmatrix

$$\underline{X} = \begin{pmatrix} 1 & x_{11} & x_{12} \\ 1 & x_{21} & x_{22} \\ \vdots & \vdots & \vdots \\ 1 & x_{10,1} & x_{10,2} \end{pmatrix} = \begin{pmatrix} 1 & 8,90 & 9,56 \\ 1 & 44,46 & 18,79 \\ 1 & 18,64 & 12,16 \\ 1 & 15,77 & 5,71 \\ 1 & 18,99 & 4,32 \\ 1 & 33,09 & 5,66 \\ 1 & 11,89 & 9,51 \\ 1 & 10,90 & 10,37 \\ 1 & 9,63 & 4,05 \\ 1 & 10,88 & 10,19 \end{pmatrix}$$

und die Vektoren

$$\underline{y} = \begin{pmatrix} 0,14 \\ 2,57 \\ 0,47 \\ 1,01 \\ 1,52 \\ 2,74 \\ 0,41 \\ 0,41 \\ 1,71 \\ 0,29 \end{pmatrix}, \underline{b} = \begin{pmatrix} b_0 \\ b_1 \\ b_2 \end{pmatrix} \text{ und } \underline{d} = \begin{pmatrix} d_1 \\ d_2 \\ \dots \\ d_{10} \end{pmatrix}.$$

Mit diesen Matrizen können wir das Gleichungssystem in Matrixschreibweise darstellen als

$$\underline{y} = \underline{X}\underline{b} + \underline{d}$$

wobei Vektoren und Matrizen unterstrichen sind.

Auch hier wird die Quadratsumme der Residuen minimiert, um die Regressionskoeffizienten zu erhalten. Diese berechnen sich mit der Formel

$$\underline{b} = (\underline{X}'\underline{X})^{-1}\underline{X}'\underline{y}.$$

Wir erhalten den Vektor der Regressionskoeffizienten

$$\underline{b} = (\underline{X}'\underline{X})^{-1}\underline{X}'\underline{y} = \begin{pmatrix} 0,556 \\ 0,089 \\ -0,117 \end{pmatrix},$$

also

$$\begin{aligned} \text{popgrow} &= 0,556 + 0,089 \cdot \text{birth} - 0,117 \cdot \text{mort} \text{ bzw.} \\ \hat{y} &= 0,556 + 0,089x_1 - 0,117x_2. \end{aligned}$$

Das Bestimmtheitsmaß ist hier 88,8%. Die Anpassung hat sich beträchtlich verbessert. Hier scheint das Zusammenwirken der beiden Regressoren mehr bewirkt zu haben als die „Summe der Teile“.

Die Wurzel aus dem Bestimmtheitsmaß ergibt den **multiplen Korrelationskoeffizienten** $r = 0,942$. Der multiple Korrelationskoeffizient kann nur zwischen 0 und 1 liegen, wobei 1 wieder vollständige Korrelation bedeutet.

Die Regressionskoeffizienten 0,089 und 0,117 sind die **partiellen Ableitungen** der Regressionsebene. Man könnte die Koeffizienten so interpretieren: Steigt bei konstanter Sterblichkeit die Geburtenrate um einen Punkt, erhöht sich das Bevölkerungswachstum um ca. 0,1 Prozent.

Steigt dagegen bei konstanter Geburtenrate die Sterblichkeit um einen Punkt, sinkt das Bevölkerungswachstum um ca. einen Punkt. Eine simultane Analyse der Regressionsebene bezüglich beider Regressionskoeffizienten ist kompliziert und meistens auch nicht sinnvoll interpretierbar. Die Analyse eines Regressionskoeffizienten bei Konstanzhaltung der übrigen Regressoren nennt man eine **Ceteris-Paribus-Analyse**.

In der Regel ist die Berechnung der Regressionskoeffizienten im multiplen linearen Regressionsmodell so aufwendig, daß Computer eingesetzt werden müssen. Spezielle statistische Datenbanksysteme wie SPSS,SAS oder Minitab ermöglichen eine umfassende Regressionsanalyse.

Die Vor- und Nachteile der Minimum-Quadrat-Methode sind dieselben wie bei der Einfachregression: Es sei $x_{8,2} = 100$ statt 10,9. Man erhält

$$popgrow = 1,13 + 0,0031 \cdot birth - 0,0092 \cdot mort$$

mit einem Bestimmtheitsmaß von 0,7%.

5.3.

Zeitreihenanalyse

[<< hoch zum Anfang vom Inhaltsverzeichnis](#)

[<< hoch zum Kapitel 5 vom Inhaltsverzeichnis](#)

Einführung

Zeitreihen sind Beobachtungen, die im Lauf der Zeit erhoben wurden. Bei der Analyse von Zeitreihen versuchen wir, die Beobachtungen durch den Faktor Zeit zu erklären. Wir suchen nach bestimmten Gesetzmäßigkeiten, nach denen diese Zeitreihen zustande kommen.

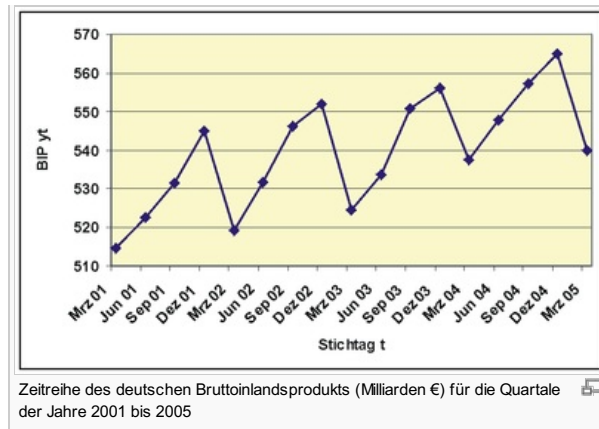
Für die optische Unterstützung stellen wir eine Zeitreihe als Streudiagramm dar. Um den Verlauf, die Entwicklung des Merkmals darstellen, können wir die Punkte zu einer Kurve (Polygonzug) verbinden.

Wir haben hier beispielsweise das Bruttoinlandsprodukt der Bundesrepublik Deutschland (Quelle: © Statistisches Bundesamt Deutschland 2005) der Quartale 2001 bis 2005 gegeben.

Stichtag	Mrz 01	Jun 01	Sep 01	Dez 01	Mrz 02	Jun 02
BIP	514,51	522,63	531,51	544,91	519,19	531,66

Stichtag	Sep 02	Dez 02	Mrz 03	Jun 03	Sep 03	Dez 03
BIP	546,06	551,9	524,4	533,59	550,76	556,12

Stichtag	Mrz 04	Jun 04	Sep 04	Dez 04	Mrz 05
BIP	537,36	547,85	557,21	564,82	539,78



Modell der Zeitreihe

Die Zeitreihenanalyse erfordert die Konzipierung der Zeitreihe als Modell:

Wir betrachten einen Beobachtungszeitraum mit T vielen Zeitpunkten t . Zu einem Zeitpunkt t gehört die Beobachtung y_t des Merkmals y .

Da Zeitangaben häufig unhandlich bei der Berechnung sind (z. B. 1.3.1996), empfiehlt es sich, die Zeitpunkte durchzunummerieren, z.B. $t = 1, 2, \dots, n$.

Beispiel Großhandel

Es liegen $n = 60$ Quartalsumsätze des Gartenbedarfsgroßhandels Rosalinde vor. Die Quartale sind durchnummeriert als $t = 1, \dots, 60$. Es sind hier nur die ersten Beobachtungen wiedergegeben. Die komplette Zeitreihe befindet sich in [Zeitreihe Rosalinde](#).

Stichtag zum Ende des Monats	Quartal	Umsatz in Mio. €	Linearer Trend
Mrz 90	1	52,19	42
Jun 90	2	48,69	44
Sep 90	3	49,28	46
...

Stichtag zum Ende des Monats	Saisonaler Zyklus	Konjunktureller Zyklus	Restschwankung
Mrz 90	6,00	3,06	1,13
Jun 90	0,00	5,66	-0,96
Sep 90	-6,00	7,39	1,89
...

Wir sehen, dass die Tendenz der Umsätze steigend ist. Es scheint

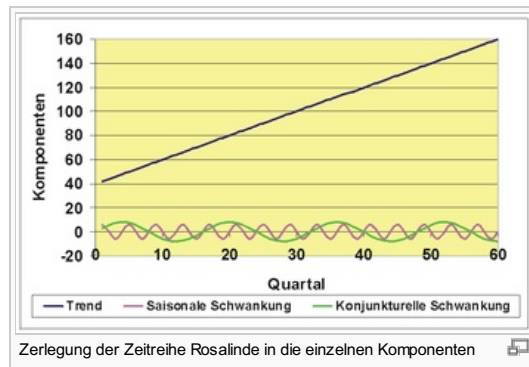
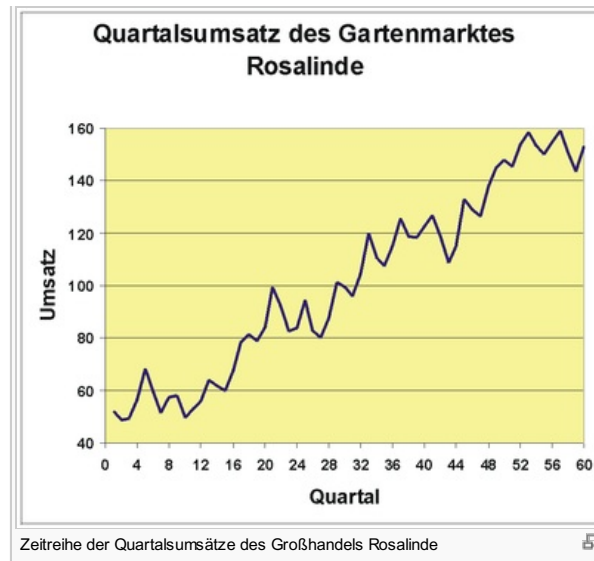
Wir sehen, dass die Tendenz der Umsätze steigend ist. Es scheint sich außerdem ein vermutlich konjunktureller Zyklus abzuzeichnen, der z. B. 1992 ein Tief und 1995 ein Hoch hatte. Und es ist deutlich ein einjähriger, saisonaler Zyklus zu erkennen, der auch aus der Tabelle ersichtlich ist.

Wir können also die Komponenten der Zeitreihe unterscheiden:

- Trend Q
- Konjunkturelle Schwankung K
- Saisonale Schwankung S
- Restschwankung r

Sind diese Komponenten unabhängig voneinander, gehen wir vom additiven Modell aus:

$$y = Q + K + S + r$$



Oft überlagern sich mehrere zyklische Schwankungen. Es gibt hier spezielle Verfahren, die Zyklen zu identifizieren.

Ein Problem in der Zeitreihenanalyse ist die Wahl des richtigen Modells. Bei den einfacheren Modellen beschränkt man sich meist auf die Bestimmung einer glatten Komponente, die aus Trend und/oder konjunktureller Komponente gebildet wird, einer saisonalen Komponente und die Restschwankung.

5.3.1. Trend und Saisonkomponente

[<< hoch zum Anfang vom Inhaltsverzeichnis](#)
[<< hoch zum Kapitel 5 vom Inhaltsverzeichnis](#)

Üblicherweise wird bei der Schätzung des Trends Q und der Saisonkomponente S so vorgegangen, dass zuerst der Trend Q bestimmt wird. Es wird dann y vom Trend bereinigt, d.h. von Beobachtungen y_t werden die Trendwerte Q_t abgezogen. Aus den resultierenden Restwerten wird dann die saisonale Komponente errechnet. Man kann auch beide Komponenten mit Hilfe der multiplen Regression auf einmal bestimmen.

Schätzung des Trends mit der Regressionsgerade

Wenn wir von einem **linear verlaufenden Trend** ausgehen können, schätzen wir ihn mit dem Regressionsmodell

$$\hat{y}_t = a + bt \text{ bzw. } y_t = a + bt + d_t \quad (t = 1, 2, \dots, T; y_t = y_1, y_2, \dots, y_T)$$

mit den Lösungen

$$\begin{aligned} b &= \frac{\sum_{t=1}^T (t - \bar{t})(y_t - \bar{y})}{\sum_{t=1}^T (t - \bar{t})^2} \\ &= \frac{\sum_{t=1}^T t \cdot y_t - T \cdot \bar{t} \cdot \bar{y}}{\sum_{t=1}^T t^2 - T \cdot \bar{t}^2} \\ &= \frac{\sum_{t=1}^T t y_t - \frac{T(T+1)}{2} \bar{y}}{\frac{1}{12}(T^3 - T)} \end{aligned}$$

und

$$\begin{aligned} a &= \bar{y} - b \cdot \bar{t} \\ &= \bar{y} - b \cdot \frac{T+1}{2} \end{aligned}$$

Die Trendwerte Q_t sind dann

$$Q_t = \hat{y}_t = a + bt.$$

Beispiel Herrenbekleidung

Die monatlichen Aufträge für die letzten 3 Jahre eines Herstellers für Herrenbekleidung (in 1000 Stück) sind durch die unten folgende Zeitreihe in der [Zeitreihe Herrenbekleidung](#) gegeben, von der ein Ausschnitt vorliegt. Die Grafik zeigt, daß offensichtlich ein steigender linearer Trend mit saisonalem Jahreszyklus vorliegt.

t	y _t	t·y _t	t ²
1	20	20	1
2	22	44	4
3	24	72	9
4	21	84	16
...
666	828	15889	16206

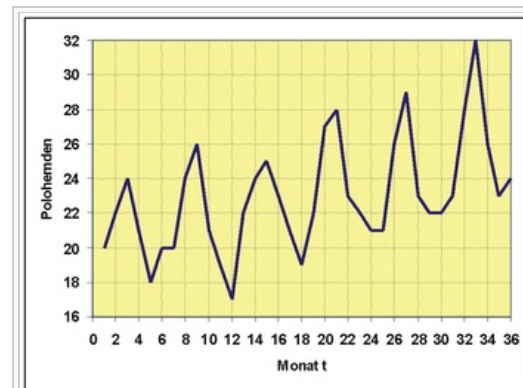
Wir ermitteln zuerst die arithmetischen Durchschnitte:

$$\bar{t} = \frac{666}{36} = 18,5 \text{ und entsprechend } \bar{y}_t = 23. \text{ Dann erhalten wir für}$$

den Regressionsansatz

$$\hat{y}_t = a + bt$$

die Regressionskoeffizienten nach dem Verschiebungssatz



Monatliche Aufträge für Polohemden eines Herstellers für Herrenbekleidung

$$b = \frac{15889 - 36 \cdot 18,5 \cdot 23}{16206 - 36 \cdot 18,5^2} = 0,1470$$

und

$$a = \bar{y} - b \cdot \bar{t} = 23 - 0,1470 \cdot 18,5 = 20,2810$$

Die geschätzten Trendwerte sind $\hat{y}_t = a + bt$, z.B.

$$\hat{y}_1 = 20,2810 + 0,1470 \cdot 1 \approx 20,43$$

$$\hat{y}_2 = 20,2810 + 0,1470 \cdot 2 \approx 20,57$$

usw.

Die Residuen sind

$$y_1 - \hat{y}_1 = 20 - 20,43 = -0,43$$

$$y_2 - \hat{y}_2 = 22 - 20,57 = 1,43$$

usw.

t	y _t	a + bt	d _t
1	20	20,43	-0,43
2	22	20,57	1,43
3	24	20,72	3,28
4	21	20,87	0,13
5	18	21,02	-3,02
6	20	21,16	-1,16
...
34	26	25,28	0,72
35	23	25,43	-2,43
36	24	25,57	-1,57

Liegt ein nichtlinearer Trendverlauf vor, kann auch ein nichtlinearer Regressionsansatz gewählt werden. Es können neben t auch andere exogene Variablen in das Modell aufgenommen werden.

Schätzung der Saisonkomponente

Gehen wir von dem additiven Modell

$$y_t = Q_t + S_t + r_t$$

aus, bleibt nach Schätzung der Trendkomponente Q noch die Abweichung

$$d_t = y_t - Q_t$$

übrig, die sich zusammensetzt aus

$$d_t = S_t + r_t$$

Wir nennen deshalb d_t auch den trendbereinigten Zeitreihenwert. Es soll nun noch die saisonale Komponente S_t ermittelt werden. Wir könnten etwa versuchen, diese zyklische Komponente mit einer Sinusfunktion zu schätzen. Einfacher ist aber folgendes Vorgehen: Wir ermitteln die trendbereinigten Zeitreihenwerte d_t. Dann wird aus allen Werten d_t, die die gleiche Saison betreffen, ein arithmetischer Durchschnitt gebildet, der als Schätzung für die saisonale Komponente verwendet wird.

Beispiel Herrenbekleidung

Für die Januar-Saisonkomponente werden alle Januarwerte der d_t gemittelt:

$$S_{jan} = S_1 = S_{13} = S_{25} = \frac{-0,43 - 0,19 - 2,96}{3} = -1,19$$

usw.

$$r_t = y_t - Q_t - S_t$$

ergibt dann die nichterklärte Restschwankung.

Wir können jetzt eine Prognose für den Zeitpunkt $T+k$ ermitteln als

$$\hat{y}_{T+k} = Q_{T+k} + S_{T+k},$$

wobei wir für S_t den Saisonwert für diejenige Saison wählen, die in $T+k$ auftritt.

Beispiel für eine Prognose:

Wir wollen für März des 4. Jahres eine Prognose des Auftragseingangs machen. Es handelt sich um den Zeitpunkt $t = 39$.

Wir erhalten den Trend als

$$Q_{39} = 20,281 + 39 \cdot 0,147 = 26,014$$

und die Saisonkomponente als

$$S_3 = \frac{3,28 + 2,51 + 4,75}{3} = 3,51.$$

Die Prognose errechnet sich nun als

$$26,014 + 3,51 = 29,524.$$

Multiplikative Verknüpfung der Zeitreihen-Komponenten

Bisher wurde von einer additiven Überlagerung des Trends durch die Saisonkomponente ausgegangen, d.h. die Komponenten wurden als unabhängig angesehen. Häufig nehmen aber die zyklischen Schwankungen mit steigendem Trend zu. Es könnte hier beispielsweise das multiplikative Modell

$$y_t = Q_t \cdot K_t \cdot r_t$$

vorliegen. Wir können den Ansatz logarithmieren und erhalten

$$\log y_t = \log Q_t + \log S_t + \log r_t$$

Mit dem logarithmierten Ansatz führen wir die Zerlegung des Modells in seine Komponenten durch, wie oben beschrieben.

5.3.2.

Glättungsverfahren

[<< hoch zum Anfang vom Inhaltsverzeichnis](#)

[<< hoch zum Kapitel 5 vom Inhaltsverzeichnis](#)

Schätzung der glatten Komponente mit gleitenden Durchschnitten

Lässt sich die Trendkomponente des Zeitreihenmodells offensichtlich durch keine funktionale lineare oder nichtlineare Beziehung darstellen, kann man eine **glatte Komponente** mit Hilfe gleitender Durchschnitte bestimmen.

Gleitende Durchschnitte ungeradzahligter Ordnung

Beispiel Hotelaufenthalte (G. D. 3. O)

In einem Kurhotel werden Ende April, Ende August und Ende Dezember die Zahl der Hotelaufenthalte festgehalten. Es wurde mit Ende Dezember begonnen.

Stichtag	t	Aufenthalte y_t
Dez 89	1	408
Apr 90	2	372
Aug 90	3	480
Dez 90	4	444

Apr 91	5	447
Aug 91	6	492
Dez 91	7	429
Apr 92	8	411
Aug 92	9	486
Dez 92	10	525
Apr 93	11	495

Zur Ermittlung des Trends wurden gleitende Durchschnitte 3. Ordnung gebildet:

$$\bar{y}_2 = \frac{1}{3} \cdot (408 + 372 + 480) = 420$$

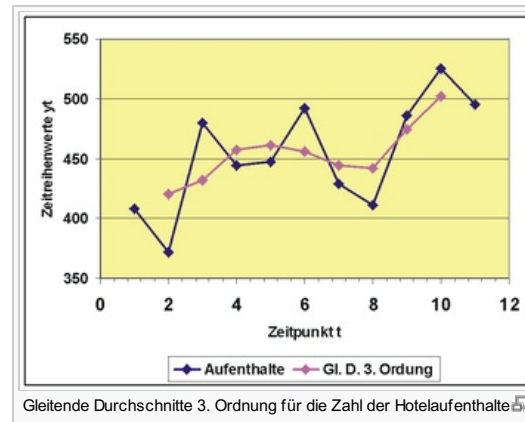
$$\bar{y}_3 = \frac{1}{3} \cdot (372 + 480 + 444) = 432$$

...

$$\bar{y}_9 = \frac{1}{3} \cdot (411 + 486 + 525) = 474$$

$$\bar{y}_{10} = \frac{1}{3} \cdot (486 + 525 + 495) = 502$$

Stichtag	t	Aufenthalte y_t	\bar{y}_t
Dez 89	1	408	
Apr 90	2	372	420
Aug 90	3	480	432
Dez 90	4	444	457
Apr 91	5	447	461
Aug 91	6	492	456
Dez 91	7	429	444
Apr 92	8	411	442
Aug 92	9	486	474
Dez 92	10	525	502
Apr 93	11	495	



Der Index t der Glättung \bar{y}_t entspricht immer dem Beobachtungswert in der Mitte der beteiligten Zeitreihenwerte.

Man sieht, dass die gleitenden Durchschnitte die starken Schwankungen glätten und man den Trend, oder besser die glatte Komponente, besser erkennt.

Die Zahl der beteiligten Beobachtungen gibt die Ordnung des Durchschnitts an. Man berechnet einen gleitenden Durchschnitt 3. Ordnung folgendermaßen:

$$\begin{aligned}\bar{y}_2 &= \frac{y_1 + y_2 + y_3}{3} \\ \bar{y}_3 &= \frac{y_2 + y_3 + y_4}{3} \\ \dots \\ \bar{y}_t &= \frac{y_{t-1} + y_t + y_{t+1}}{3} \\ \dots \\ \bar{y}_{n-1} &= \frac{y_{n-2} + y_{n-1} + y_n}{3}\end{aligned}$$

Entsprechend ergeben sich gleitende Durchschnitte 5. Ordnung als

$$\begin{aligned}\bar{y}_3 &= \frac{y_1 + y_2 + y_3 + y_4 + y_5}{5} \\ \bar{y}_4 &= \frac{y_2 + y_3 + y_4 + y_5 + y_6}{5} \\ \dots \\ \bar{y}_t &= \frac{y_{t-2} + y_{t-1} + y_t + y_{t+1} + y_{t+2}}{5} \\ \dots\end{aligned}$$

usw.

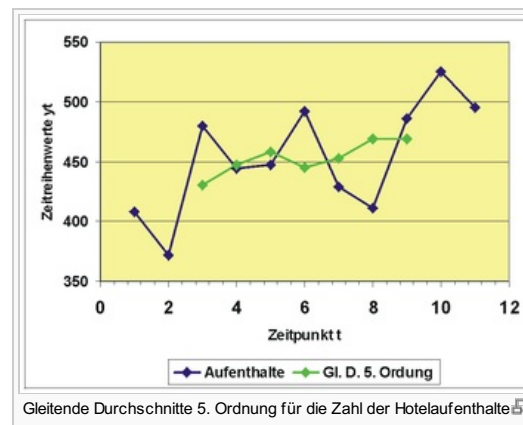
Beispiel Hotelaufenthalte (G. D. 5. O)

Wir berechnen die gleitenden Durchschnitte 5. Ordnung als

$$\begin{aligned}\bar{y}_3 &= \frac{408 + 372 + 480 + 444 + 447}{5} \\ \bar{y}_4 &= \frac{372 + 480 + 444 + 447 + 492}{5}\end{aligned}$$

usw., also

Stichtag	t	Aufenthalte y_t	\bar{y}_t
Dez 89	1	408	
Apr 90	2	372	
Aug 90	3	480	430,2
Dez 90	4	444	447
Apr 91	5	447	458,4
Aug 91	6	492	444,6



Dez 91	7	429	453
Apr 92	8	411	468,6
Aug 92	9	486	469,2
Dez 92	10	525	
Apr 93	11	495	

Zur Prognose über den Beobachtungszeitraum hinaus sind gleitende Durchschnitte nicht so recht geeignet, da die Randwerte der Zeitreihe nicht geschätzt werden. Allerdings gibt es Verfahren, mit denen man diese Werte durch eine Gewichtung der benachbarten Werte ausfüllen kann.

Gleitende Durchschnitte geradzahliger Ordnung

Die Rechtfertigung für gleitende Durchschnitte als Schätzung für die glatte Komponente begründet sich darin, daß sich saisonale Schwankungen kompensieren, d.h. daß sich positive und negative Abweichungen der Saisonkomponente vom Trend aufheben. Das ist allerdings nur dann gegeben, wenn die Breite, also die Ordnung des gleitenden Durchschnitts genau einen Zyklus umfaßt und wenn alle Zyklen die gleiche Länge aufweisen. Sonst können gleitende Durchschnitte verzerrte Schätzungen liefern.

Die Breite saisonaler Schwankungen ist allerdings meist geradzahlig, z.B. Quartale, Monate usw. Es entsteht hier das Problem, daß bei der Durchschnittsbildung ein Schätzwert $\bar{y}_{t+0,5}$ zwischen zwei Beobachtungswerte y_t und y_{t+1} zu liegen kommt, was meist unerwünscht ist. Man behilft sich hier so, daß man den Durchschnitt auf der Zeitachse um eine halbe Zeiteinheit nach rechts verschiebt. Die beiden Randwerte, die nur teilweise erfaßt werden, gehen mit einem Gewicht von 0,5 in den Durchschnitt ein,

z.B. statt

$$\bar{y}_{3,5} = \frac{y_2 + y_3 + y_4 + y_5}{4}$$

berechnet man

$$\bar{y}_4 = \frac{\frac{y_2}{2} + y_3 + y_4 + y_5 + \frac{y_6}{2}}{4}$$

Beispiel

Es liegen 12 Quartalswerte vor und es sollen gleitende Durchschnitte 4. Ordnung ermittelt werden.

t	yt	gleitender Durchschnitt 4. Ordnung
1	12	
2	8	
3	12	12
4	14	13
5	16	14,125
6	12	15,625
7	17	17,5
8	21	19
9	24	20,5

10	16	21,5
11	25	
12	21	

Z.B.

$$\bar{y}_3 = \frac{\frac{12}{2} + 8 + 12 + 14 + \frac{16}{2}}{4} = 12$$

$$\bar{y}_4 = \frac{\frac{8}{2} + 12 + 14 + 16 + \frac{12}{2}}{4} = 13$$

Exponentielle Glättung

Lässt eine Zeitreihe keinerlei systematisches Muster wie [linearen](#) Anstieg oder Ähnliches erkennen, kann man versuchen, mit der exponentiellen Glättung eine glatte Komponente nachzubilden. Insbesondere kann man damit eine Prognose für den Zeitpunkt $T + 1$ erhalten.

Das Verfahren wird beispielsweise in der [Lagerhaltung](#) verwendet, wenn es etwa darum geht, den Bedarf eines zu bestellenden Artikels im kommenden Jahr zu ermitteln. So hat etwa die [Schweizer Armee](#) mit der exponentiellen Glättung gute Erfolge bei der Ermittlung der benötigten Gewehre im folgenden Jahr gemacht.

Man geht von dem Ansatz aus, dass der gegenwärtige Zeitreihenwert immer auch von den vergangenen Werten beeinflusst wird, wobei sich der Einfluss abschwächt, je weiter der Wert in der Vergangenheit liegt.

Formales Modell

Gegeben ist eine Zeitreihe mit den Beobachtungen $y_1, y_2, \dots, y_t, \dots$ zu den Zeitpunkten t . Im Zeitpunkt t wird für y_t ein geglätteter Schätzwert y_t^* errechnet, der sich als gewichteter Durchschnitt ergibt aus dem aktuellen Zeitreihenwert y_t und dem Schätzwert der Vorperiode y_{t-1}^* . Die

Gewichtung wird durch den Glättungsfaktor α bestimmt, wobei $0 \leq \alpha \leq 1$ sein muss. Man erhält

$$y_t^* = \alpha \cdot y_t + (1 - \alpha) \cdot y_{t-1}^*.$$

Die Zeitreihe baut sich so rekursiv auf. Theoretisch ist die laufende Zeitreihe beim Zeitpunkt t bereits unendlich lang. Für die praktische Ermittlung des geglätteten Wertes wird man allerdings einen Startwert y_0^* vorgeben und von da an die geglättete Zeitreihe ermitteln.

Baut man nun, beginnend bei y_0^* , die geglättete Zeitreihe auf,

$$\begin{aligned} y_1^* &= \alpha y_1 + (1 - \alpha) y_0^*, \\ y_2^* &= \alpha y_2 + (1 - \alpha) y_1^*, \\ y_3^* &= \alpha y_3 + (1 - \alpha) y_2^*, \\ &\dots \end{aligned}$$

erhält man, wenn man die Rekursivität auflöst,

$$y_t^* = \alpha y_t + \alpha(1 - \alpha)y_{t-1} + \alpha(1 - \alpha)^2 y_{t-2} + \dots + \alpha(1 - \alpha)^{t-1} y_1 + (1 - \alpha)^t y_0.$$

Man sieht, wie wegen $\alpha < 1$ die Einflüsse der Vergangenheit immer mehr verschwinden.

Der Schätzwert y_t^* liefert den Prognosewert für den Zeitpunkt $t+1$. Liegt dann im Zeitpunkt $t + 1$ eine neue Beobachtung vor, kann die Prognose für $t + 2$ ermittelt werden usw.

Für die Wahl des Glättungsfaktors wird häufig 0,2 bis 0,3 empfohlen. Man kann aber auch mit Hilfe der Regressionsanalyse den Glättungsfaktor schätzen.

Einfaches Zahlenbeispiel

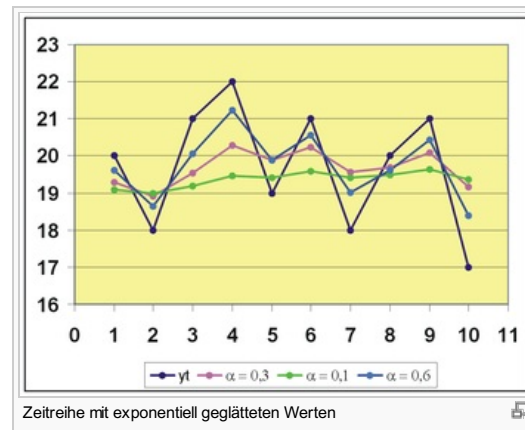
Es sind die Zeitreihenwerte y_1, \dots, y_{10} gegeben, wie unten in der Tabelle aufgeführt. Diese Werte sollen exponentiell geglättet werden. Es wurde ein Glättungskoeffizient von $\alpha = 0,3$ gewählt und man benötigt einen Anfangswert, der hier $y_0^* = 19$ betragen soll. Wir beginnen

$$y_1^* = 0,3 \cdot 20 + 0,7 \cdot 19 = 6 + 13,3 = 19,3$$

$$y_2^* = 0,3 \cdot 18 + 0,7 \cdot 19,3 = 18,91$$

usw. In der Tabelle sind die Glättungen für ausgewählte Werte von α aufgeführt.

t	y_t	y_t^*	y_t^*	y_t^*
		$\alpha = 0,3$	$\alpha = 0,1$	$\alpha = 0,6$
0	-	19	19	19
1	20	19,3	19,1	19,6
2	18	18,91	18,99	18,64
3	21	19,54	19,19	20,06
4	22	20,28	19,47	21,22
5	19	19,89	19,42	19,89
6	21	20,23	19,58	20,56
7	18	19,56	19,42	19,02
8	20	19,69	19,48	19,61
9	21	20,08	19,63	20,44
10	17	19,16	19,37	18,38

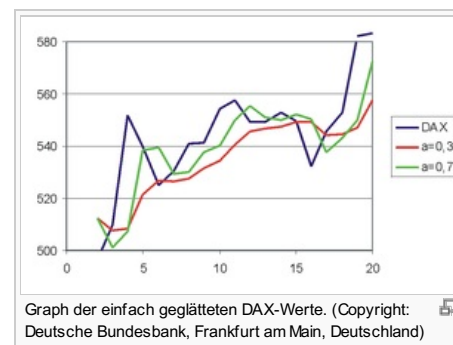


Die Graphik zeigt die Glättung für $\alpha = 0,1$, $\alpha = 0,3$ und $\alpha = 0,6$. Man sieht, dass der kleinste Glättungsfaktor die Zeitreihe am stärksten glättet, denn hier geht der aktuelle Wert nur mit einem Gewicht von 0,1 ein, wogegen die „mittleren“ Vergangenheitswerte weiterhin mit 0,9 berücksichtigt werden.

Beispiel für den exponentiell geglätteten DAX

Es soll mit den monatlichen Durchschnittswerten des Aktienindex DAX für die Monate Januar 1977 bis August 1978 eine exponentielle Glättung berechnet werden. Die Daten liegen nebst den geglätteten Zeitreihenwerten in der Tabelle vor:

DAX-Werte und ihre exponentielle Glättung ($\alpha = 0,3$)			
Monat	Zeitpunkt t	DAX y_t	Glättung y_t^*
1977 Jan	0	512,3	512,3
1977 Feb	1	496,2	507,5
1977 Mrz	2	509,8	508,2
1977 Apr	3	551,9	521,3
1977 Mai	4	539,9	526,9



1977 Jun	5	524,9	526,3
1977 Jul	6	530,3	527,5
1977 Aug	7	540,9	531,5
1977 Sep	8	541,3	534,4
1977 Okt	9	554,2	540,4
1977 Nov	10	557,5	545,5
1977 Dez	11	549,34	546,7
1978 Jan	12	549,4	547,5
1978 Feb	13	552,9	549,1
1978 Mrz	14	549,7	549,3
1978 Apr	15	532,1	544,1
1978 Mai	16	545,5	544,5
1978 Jun	17	553,0	547,1
1978 Jul	18	582,1	557,6
1978 Aug	19	583,1	565,2

Der erste Wert wird mit 512,3 als Startwert y^*_0 genommen. Wir verwenden einen Glättungsfaktor $\alpha = 0,3$.

Es ergeben sich die geglätteten Werte

$$\begin{aligned} y_1^* &= 0,3 \cdot 496,2 + 0,7 \cdot 512,3 = 507,5, \\ y_2^* &= 0,3 \cdot 509,8 + 0,7 \cdot 507,5 = 508,2, \\ y_3^* &= 0,3 \cdot 551,9 + 0,7 \cdot 508,2 = 521,3, \\ &\dots \end{aligned}$$

Die Schätzung y^*_1 ist jetzt der Prognosewert für die Periode 2 und so weiter.

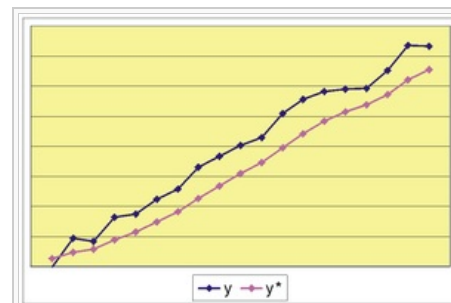
Die Grafik zeigt die Glättung für $\alpha = 0,3$ und $\alpha = 0,7$. Man sieht, dass der kleinere Glättungsfaktor die Zeitreihe stärker glättet, denn hier geht der aktuelle Wert jetzt nur mit einem Gewicht von 0,3 ein, wogegen die „mittleren“ Vergangenheitswerte weiterhin mit 0,7 berücksichtigt werden.

Exponentielle Glättung bei trendbehafteten Werten

Die exponentielle Glättung ist dann ein empfehlenswertes Verfahren, wenn die Zeitreihenwerte einen chaotischen Eindruck machen und keinerlei Systematik erkennen lassen. Liegen allerdings Beobachtungen vor, die einen Trend beinhalten, d.h. die laufend steigen oder fallen, „schleppen“ die geglätteten Werte „hinterher“. Man sieht in der Grafik deutlich, wie die Schätzwerte immer systematisch unter den beobachteten Werten liegen. Auch in der oberen Grafik der Dax-Kurse liegen beispielsweise zwischen $t = 7$ und $t = 12$ die einfach geglätteten Schätzwerte immer systematisch unter den beobachteten Werten.

Eine zufriedenstellende Lösung für das Problem, daß bei einem steigenden (fallenden) Trend die Zeitreihenwerte systematisch unterschätzt (überschätzt) werden, bieten gleitende Durchschnitte zweiter Ordnung. Hier werden die bereits einmal geglätteten Werte noch einmal einer Glättung unterzogen. Man erhält den Schätzwert y^{**} , der sich analog zu oben berechnet aus

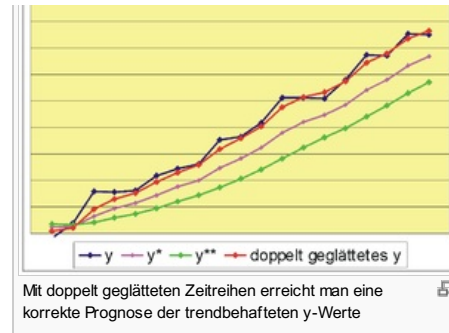
$$y_t^{**} = \alpha \cdot y_t^* + (1 - \alpha) \cdot y_{t-1}^{**}$$



Die geglätteten Prognosewerte y^* liegen systematisch unter den beobachteten trendbehafteten Zeitreihenwerten y

Für einen brauchbaren Prognosewert für Periode t+1 muss man dann bestimmen

$$\hat{y}_{t+1} = 2 \cdot y_t^* - y_{t-1}^{**}$$



6.	Maßzahlen	<< hoch zum Anfang vom Inhaltsverzeichnis << hoch zum Kapitel 6 vom Inhaltsverzeichnis
6.1.	Konzentration	<< hoch zum Anfang vom Inhaltsverzeichnis << hoch zum Kapitel 6 vom Inhaltsverzeichnis

Die **Konzentration** befasst sich mit der Intensität, mit der sich ein Objekt auf eine vorgegebene Menge verteilt. Eine typische Aussage der Konzentrationsmessung wäre etwa: 20% der Menschen eines bestimmten Staates besitzen 90% des Vermögens. Demnach teilen sich die anderen 80% die restlichen 10%. Hier kann man von einer starken Konzentration sprechen.

Kino-Beispiel

Im Rahmen einer Controllinganalyse eines Kinos wurden die Besucherzahlen (Merkmal x) für die 5 angebotenen Spielfilme an einem Tag erfasst. Man erhielt die Tabelle

Filmtitel	Zahl der Besucher x
Rotkäppchen	25
Verliebt ins Abendrot	75
Leif Erikson	125
Söhne der Alhambra	250
Galaxy-Fighter	525

Definitionen

Es gibt verschiedene Verfahren zur Konzentrationsmessung. Man kann die Konzentration grafisch darstellen oder Kennwerte berechnen. Die Merkmalsbeträge x müssen aufsteigend geordnet vorliegen, also $x_{[1]} \leq x_{[2]} \leq \dots \leq x_{[n]}$.

Für die Konzentrationsmessung werden neben der relativen Summenfunktion S_i^* folgende Definitionen benötigt:

- Merkmalssumme $\sum_i x_i = n \cdot \bar{x}$
- Kumulierte Merkmalsbeträge $q_i = \sum_{k=1}^i x_{[k]}$
- Relative kumulierte Merkmalsbeträge $q_i^* = \frac{q_i}{n \cdot \bar{x}}$

Grafik

Die Lorenzkurve ist eine grafische Darstellung der Konzentration:

Die Wertepaare $(S_i^*; q_i^*)$ werden in einem Diagramm abgetragen. Das erste Wertepaar ist $(0;0)$, das letzte $(1;1)$. Es wird zwischen diesen beiden Wertepaaren die Winkelhalbierende des Koordinatensystems eingetragen. Alle Wertepaare $(0;0)$, $(S_1^*; q_1^*)$, \dots , $(1;1)$ werden geradlinig verbunden.

Tabelle Die für die Lorenzkurve benötigten Zwischenwerte werden in der folgenden Tabelle aufgeführt. So ergibt sich beispielsweise für die kumulierten Merkmalsbeträge q_i

$$q_1 = 25, q_2 = 25 + 75 = 100, q_3 = 100 + 125 = 225 \text{ usw.}$$

Die relativen oder anteiligen Merkmalsbeträge errechnen sich durch Teilen des Gesamtmerkmalbetrags 1000, also

$$q_1^* = \frac{25}{1000} = 0,025 \text{ usw.}$$

Ebenso ermitteln wir die absolute Summenhäufigkeiten als Zahl der Filme, also

$$S_1 = 1, S_2 = 1 + 1 = 2, S_3 = 2 + 1 = 3 \dots$$

und wiederum die relative Summenhäufigkeit mit

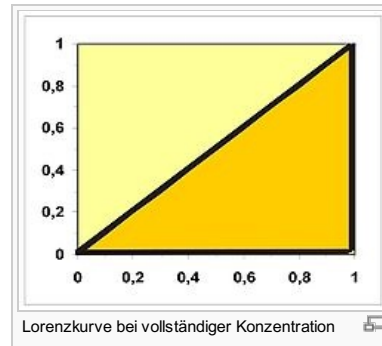
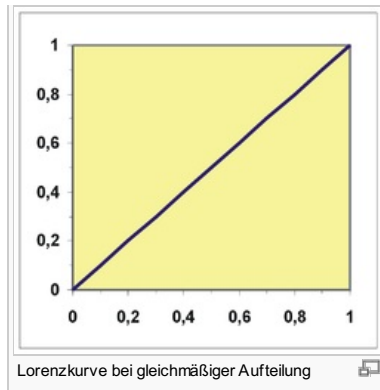
$$S_1^* = \frac{1}{5} = 0,2, S_2^* = \frac{2}{5} = 0,4, \dots$$

Es wurde außerdem noch als Platzhalter die Zeile für $i = 0$ eingefügt.

i	Filmtitel	x_i	q_i	q_i^*	S_i	S_i^*
0		0	0	0	0	0
1	Rotkäppchen	25	25	0,025	1	0,2
2	Verliebt ins Abendrot	75	100	0,100	2	0,4
3	Leif Erikson	125	225	0,225	3	0,6
4	Söhne der Alhambra	250	475	0,475	4	0,8
5	Galaxy-Fighter	525	1000	1,000	5	1
Summe		1000				

So wurden beispielsweise 40% (S_2^*) der Filme von nur 10% (q_2^*) der Besucher angesehen.

Die Lorenzkurve ist ein grafisches Maß für das Ausmaß einer Konzentration. Je weiter die Kurve „durchhängt“, desto größer ist die Konzentration. Unten sind die beiden extremen Situationen dargestellt, die gleichmäßige Aufteilung der Objekte auf die gesamte Menge und die vollständige Konzentration, bei der ein Element alle Objekte auf sich vereint und alle anderen Elemente leer ausgehen.



Werden mehrere gleichartige Gesamtheiten gegenübergestellt, bieten die verschiedenen Lorenzkurven eine schnelle optische Vergleichsmöglichkeit. Siehe dazu auch das weiter unten folgende Beispiel mit den Agrarflächen in Bayern.

Ginikoeffizient

Als Ginikoeffizient G wird bezeichnet der Anteil der Fläche, die durch die Winkelhalbierende und die Lorenzkurve gebildet wird, an der Gesamtfläche unter der Winkelhalbierenden. Wenn vollkommene Konzentration besteht, ist die Fläche über der Lorenzkurve deckungsgleich mit dem Dreieck unter der Winkelhalbierenden. G ist dann 1. Bei fehlender Konzentration ist dann $G=0$.

Ermittlung des Ginikoeffizienten

Verbindet man die Punkte auf der Lorenzkurve mit den entsprechenden Punkten auf der Winkelhalbierenden, wird klar, dass wir es mit n vielen Trapezen zu tun haben, deren Flächen wir einzeln bestimmen und dann aufsummieren. Die Fläche eines Trapezes, wie in der Grafik angegeben, ermittelt man als

$$F = \frac{1}{2} \cdot (a + c) \cdot h$$

Wir wollen die Fläche F_3 des Trapezes zwischen den Abszissenwerten (x-Achse) 0,4 und 0,6 ermitteln. Man sieht, dass das Trapez im Vergleich zur obigen Grafik gekippt vorliegt. Die Höhe h ist also die Differenz

$$S_3 - S_2 = 0,6 - 0,4 = 0,2.$$

Wir fassen a als linke Senkrechte von F_3 als a auf. Dann ist

$$a = 0,4 - 0,1 = 0,3.$$

Entsprechend beträgt die rechte Seite c

$$c = 0,6 - 0,225 = 0,375$$

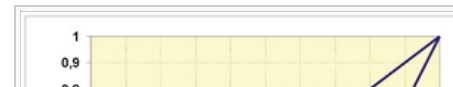
und wir erhalten als Fläche

$$F_2 = (0,3 + 0,375) \cdot 0,5 \cdot 0,2 = 0,0675.$$

Allgemein: Die obige Fläche ergibt sich dann als

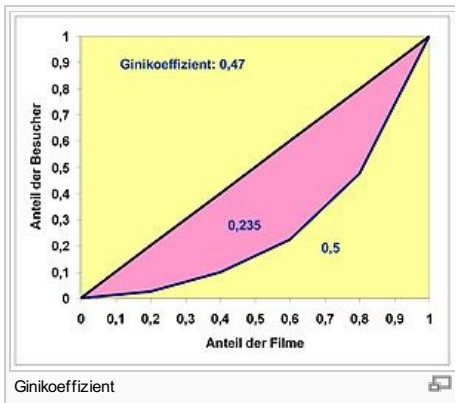
$$\sum_{i=1}^n (S_i - S_{i-1}) \cdot \frac{1}{2} ((S_i - q_i) + (S_{i-1} - q_{i-1}))$$

Es folgt beispielhaft die Berechnung des Gini in der Tabelle. Mit Tabellenkalkulation kann der Ginikoeffizient leicht ermittelt werden. Wir erhalten schließlich für den Ginikoeffizienten



$$G = \frac{0,235}{0,5} = 0,47$$

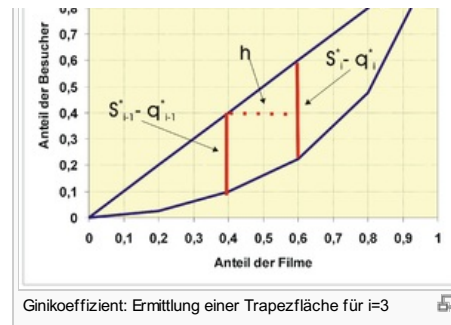
i	q_i^*	S_i^*	$h_i^* = S_i^* - S_{i-1}^*$	$a_i = S_i^* \cdot q_i^*$	$c_i = S_{i-1}^* - q_{i-1}^*$	$0,5 \cdot (a_i + c_i)$	$0,5 \cdot (a_i + c_i) \cdot h_i$
--	0	0	-	-	-	-	-
1	0,025	0,2	0,2	0,175	0	0,0875	0,0175
2	0,1	0,4	0,2	0,3	0,175	0,2375	0,0475
3	0,225	0,6	0,2	0,375	0,3	0,3375	0,0675
4	0,475	0,8	0,2	0,325	0,375	0,35	0,07
5	1	1	0,2	0	0,325	0,1625	0,0325
Summe							0,235



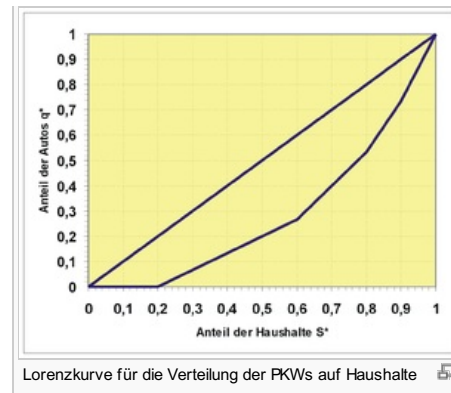
Metrisches Merkmal mit wenig möglichen Ausprägungen

Beispiel

Das interessierende Merkmal ist die Zahl der Autos in einem Haushalt. Es wurden 50 Haushalte befragt.



j	x_j	n_j	S_j	S_j^*	$x_j n_j$	q_j	q_j^*
1	0	10	10	0,2	0	0	0,00
2	1	20	30	0,6	20	20	0,27
3	2	10	40	0,4	20	40	0,53
4	3	5	45	0,9	15	55	0,73
5	4	5	50	1	20	75	1
Summe		50			75		



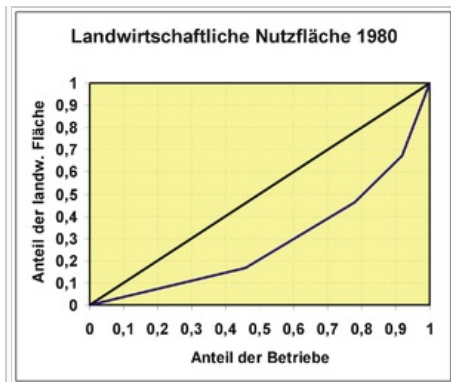
Lorenzkurve und der Ginikoeffizient berechnen sich im Prinzip wie oben, statt i wird hier der Index j verwendet. Der Merkmalsbetrag x_i wird durch $x_j^* n_j$ ersetzt.

Klassiertes Merkmal

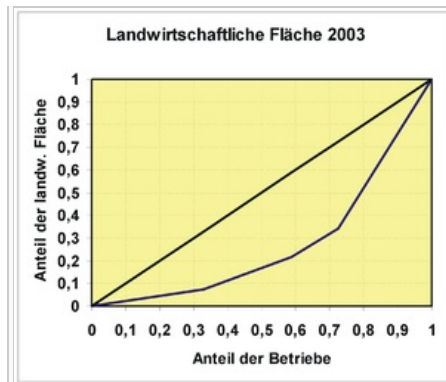
Hier wird die Klassenmitte x_j^* als Ersatz für den Merkmalswert x_j verwendet.

Beispiel

Landwirtschaftliche Nutzfläche	Zahl der Betriebe (1000)	
	1980	2003
von ... bis ... unter		
2 - 10	112	43
10 - 20	78	34
20 - 30	34	18
30 oder mehr	20	36



Lorenzkurve der Nutzfläche eines bayerischen Landwirtschaftsbetriebes im Jahr 1980



Lorenzkurve der Nutzfläche eines bayerischen Landwirtschaftsbetriebes im Jahr 2003

Klasse j von ... bis unter ...	Klassen- mitte x_j	n_j	$x_j \cdot n_j$	S_j	S_j^*	q_j	q_j^*
2 - 10	6	112	672	112	0,4590	672	0,1683
10 - 20	15	78	1170	190	0,7787	1842	0,4614
20 - 30	25	34	850	224	0,9180	2692	0,6743
30 - 100	65	20	1300	244	1,0000	3992	1,0000
Summe		244	3992				

Wir erhalten als Ginikoeffizient für das Jahr 1980 den Wert 0,43 und für das Jahr 2003 den Wert 0,46.

7.	Schätzen und Testen	<< hoch zum Anfang vom Inhaltsverzeichnis << hoch zum Kapitel 7 vom Inhaltsverzeichnis
-----------	----------------------------	---

7.1.	Konfidenzintervalle	<< hoch zum Anfang vom Inhaltsverzeichnis << hoch zum Kapitel 7 vom Inhaltsverzeichnis
-------------	----------------------------	---

Der frühere Inhalt wurde in [Mathematik: Statistik: Prinzip des Konfidenzintervalls](#) und [Mathematik: Statistik: Ausgewählte Konfidenzintervalle](#) aufgegliedert. Hier wird demnächst etwas Allgemeines über Konfidenzintervalle stehen.

7.1.1.	Prinzip des Konfidenzintervalls	<< hoch zum Anfang vom Inhaltsverzeichnis << hoch zum Kapitel 7 vom Inhaltsverzeichnis
---------------	--	---

Beispiel mit Absatz von Kaffeepaketen

Beispiel:

Eine Kaffeerösterei möchte eine neue Röstanlage anschaffen. Um beurteilen zu können, ob die Firma den aufzunehmenden Kredit tilgen kann, braucht sie Informationen über den durchschnittlichen monatlichen Absatz an Kaffeepaketen. Pro Monat muss die Firma 20.000 € Annuität zahlen. Zusammen mit den Produktionskosten sollte sie im Durchschnitt auf einen Absatz von 100.000 Kaffeepaketen im Monat kommen. Die Frage ist nun, wird sie es schaffen?

Plausible Überlegungen zur Schätzung

Der durchschnittliche monatliche Absatz von Kaffeepaketen ist unbekannt. Wie könnte man den Durchschnitt ermitteln? Man könnte eine Stichprobe mit z.B. $n = 50$ Beobachtungen ziehen und versuchen, aus dem arithmetischen Mittel \bar{x} auf den durchschnittlichen monatlichen Absatz der Grundgesamtheit zu schließen. Ist die Stichprobe groß genug, kann man vermuten, dass der Durchschnitt \bar{x} in der Grundgesamtheit, hier μ , in der Nähe von \bar{x} liegen müsste. Meistens wird \bar{x} in der Nähe von μ liegen, da aber \bar{x} die Realisation einer Zufallsvariablen ist, kann in sehr wenigen Fällen \bar{x} auch extrem weit von μ weg liegen, so daß man dann μ verkehrt einschätzt.

Wir betrachten nun den monatlichen Absatz von Kaffeepaketen (in 1000). Wir bezeichnen ihn als Zufallsvariable X . Es soll der monatliche durchschnittliche Absatz der Kaffeepäckchen geschätzt werden. Bekannt ist lediglich, dass die Zahl der verkauften Kaffeepakete normalverteilt ist mit einer Varianz 200 [1000² Stück²].

Wie sollen wir nun μ eingrenzen? Wir könnten etwa ein Intervall bestimmen, in dem z.B. 95% aller möglichen \bar{x} -Werte liegen, also

$$P(\bar{x}_u \leq \bar{X} \leq \bar{x}_o) = 0,95.$$

Damit man dieses Intervall berechnen kann, müssen Informationen über die Verteilung von \bar{X} verfügbar sein. Es soll eine Stichprobe von $n = 50$ gezogen werden, d.h. es werden die verkauften Kaffeepakete der letzten 50 Monate erfasst:

$$\bar{X} = \frac{1}{50} \sum_{i=1}^{50} X_i.$$

Verteilung des Merkmals und der Schätzfunktion

Die Zufallsvariable X in der Grundgesamtheit soll normalverteilt sein mit dem Durchschnitt $EX = \mu$ und der Varianz $\text{Var}X = \sigma^2$. Die Varianz soll bekannt sein.

Es wird eine Stichprobe vom Umfang n gezogen. Der Stichprobendurchschnitt \bar{X} ist selbst eine Zufallsvariable und ist als lineare Transformation von X wiederum normalverteilt, und zwar mit den Parametern

$$E\bar{X} = \mu \text{ und } \text{var}\bar{X} = \frac{\sigma^2}{n}.$$

Hier ist

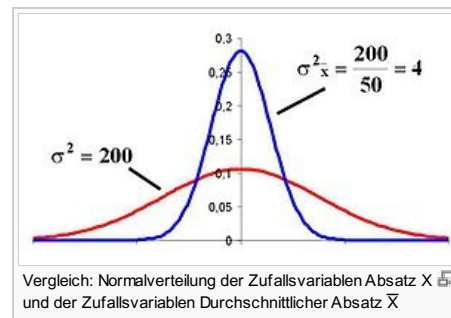
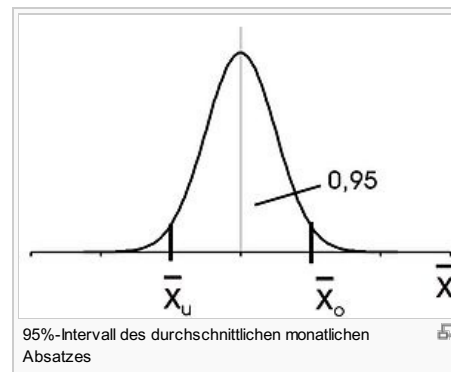
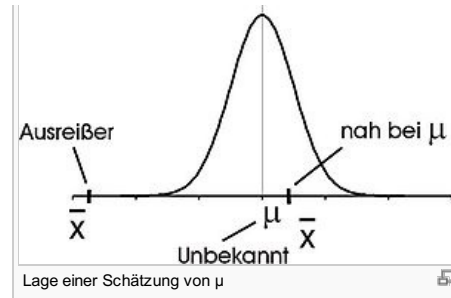
$$\frac{\sigma^2}{n} = \frac{200}{50} = 4.$$

Herleitung des Intervalls

Ausgegangen wird von

$$P(\bar{x}_u \leq \bar{X} \leq \bar{x}_o) = 0,95.$$

Untergrenze u und Obergrenze o sollen nun bestimmt werden. Wir standardisieren zunächst



$$Z = \frac{\bar{X} - \mu}{\sqrt{\frac{\sigma^2}{n}}} = \frac{\bar{X} - \mu}{\sqrt{4}},$$

so dass sich analog zu oben

$$P(z_u \leq Z \leq z_o) = 0,95$$

ergibt. z_o ist hier das 0,975-Quantil der Standardnormalverteilung. Ein Blick in die

Normalverteilungstabelle verrät uns, dass der z-Wert, der zur Wahrscheinlichkeit 0,975 gehört, 1,96 ist.

Wir können jetzt das entsprechende Intervall für Z

$$P(-1,96 \leq Z \leq 1,96) = 0,95$$

angeben. Die Ungleichung wird bezüglich μ aufgelöst:

$$P\left(-1,96 \leq \frac{\bar{X} - \mu}{\sqrt{\frac{\sigma^2}{n}}} \leq 1,96\right) = 0,95.$$

$$P\left(-1,96 \leq \frac{\bar{X} - \mu}{\sqrt{4}} \leq 1,96\right) = 0,95.$$

$$P(-1,96 \cdot 2 \leq \bar{X} - \mu \leq 1,96 \cdot 2) = 0,95.$$

$$P(-\bar{X} - 1,96 \cdot 2 \leq -\mu \leq -\bar{X} + 1,96 \cdot 2) = 0,95.$$

$$P(\bar{X} + 1,96 \cdot 2 \geq \mu \geq \bar{X} - 1,96 \cdot 2) = 0,95.$$

$$P(\bar{X} - 1,96 \cdot 2 \leq \mu \leq \bar{X} + 1,96 \cdot 2) = 0,95.$$

Dieses Intervall wird Zufallsintervall genannt, weil es von einer Zufallsvariablen (\bar{X}) gebildet wird. Wir schreiben jetzt dieses Intervall mit Symbolen:

$$P\left(\bar{X} - z \left(1 - \frac{\alpha}{2}\right) \cdot \frac{\sigma}{\sqrt{n}} \leq \mu \leq \bar{X} + z \left(1 - \frac{\alpha}{2}\right) \cdot \frac{\sigma}{\sqrt{n}}\right) = 1 - \alpha.$$

Wir bezeichnen $1 - \alpha = 0,95$ als Konfidenzkoeffizient. $\alpha = 0,05$ dagegen ist die Irrtumswahrscheinlichkeit oder das Signifikanzniveau.

Die Breite des Intervalls ist hier

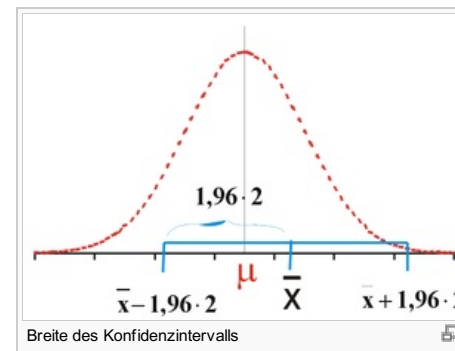
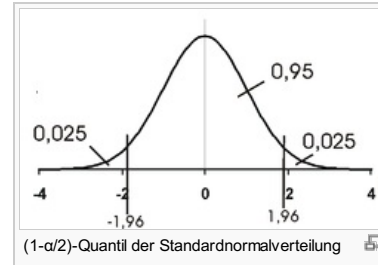
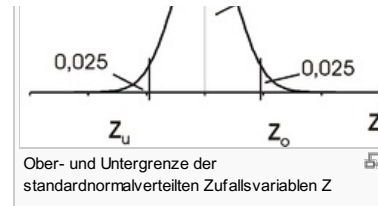
$$2 \cdot (2 \cdot 1,96) = 2 \cdot \frac{\sigma}{\sqrt{n}} \cdot z(0,975) = 7,84.$$

Also schwankt eine \bar{X} -Schätzung für μ mit einer 95%igen Wahrscheinlichkeit in einem Intervall der Breite von 7840 Kaffeepaketen, d.h. μ befindet sich mit einer 95%igen Wahrscheinlichkeit in diesem Intervall.

Es kann aber passieren, dass die Schätzung \bar{x} extrem daneben liegt. In der Grafik wurde mit \bar{x} daneben gegriffen; dieser Fall durch diese restlichen 5% abgedeckt.

Konkretes 95%-Konfidenzintervall

Es liegt nun ein konkreter Schätzwert von $\bar{x} = 98$ vor. Wir erhalten das Konfidenzintervall



$$\begin{aligned}
 & [\bar{x} - 2 \cdot 1,96; \bar{x} + 2 \cdot 1,96] \\
 &= [98 - 2 \cdot 1,96; 98 + 2 \cdot 1,96] \\
 &= [98 - 3,92; 98 + 3,92] \\
 &= [94,08; 101,92] .
 \end{aligned}$$

Entscheidung: μ kann bei einer Wahrscheinlichkeit von 95% unter 100 liegen, also kann der mittlere Umsatz unter 100.000 liegen. Deshalb sollte die Firma von dieser Investition absehen.

Was wäre, wenn man [101;108,84] erhalten hätte? Dann wäre eine dauerhafte Liquidität zu vermuten.

Einfluss der Varianz auf das Konfidenzintervall

Was wäre, wenn σ^2 statt 200 den Wert 5000 hätte? Dann wäre

$$\bar{X} \rightarrow N(\mu; \frac{5000}{50} = 100) .$$

Wir erhielten das Konfidenzintervall

$$\begin{aligned}
 & [\bar{x} - 1,96 \cdot \sqrt{100}; \bar{x} + 1,96 \cdot \sqrt{100}] \\
 &= [98 - 19,6; 98 + 19,6] \\
 &= [78,4; 117,6] .
 \end{aligned}$$

Das hieße, der wahre durchschnittliche Absatz läge mit einer Wahrscheinlichkeit von 95% zwischen 78 400 und 117 600 Päckchen. Dieses Intervall wäre eine sehr grobe Abschätzung. Mit so etwas kann man nicht mehr vernünftig planen.

Also wird das Konfidenzintervall mit steigender Varianz breiter, die Schätzungen werden schlechter, ungenauer. Hier könnte man als Abhilfe den Stichprobenumfang erhöhen.

Mindest erforderlicher Stichprobenumfang

Wie groß muss die Stichprobe mindestens sein, damit die Breite des Konfidenzintervalls höchstens 10 ist?

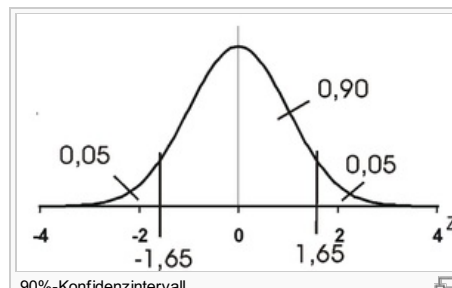
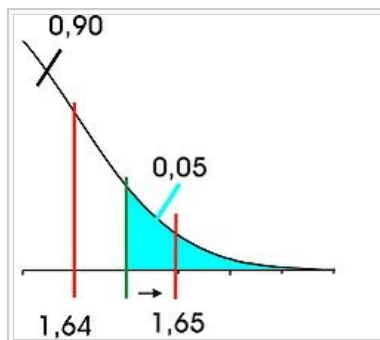
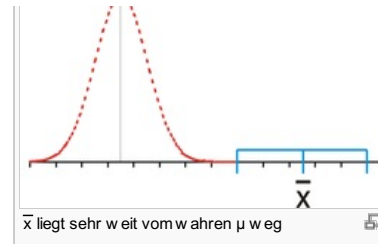
Die Breite des Konfidenzintervalls ist

$$\begin{aligned}
 2 \cdot 1,96 \cdot \sqrt{\frac{5000}{n}} &\leq 10 \rightarrow \\
 \sqrt{n} &\geq \frac{2 \cdot 1,96 \cdot \sqrt{5000}}{10} = 27,71 .
 \end{aligned}$$

Man müsste also mindestens $n = 769$ Monate erheben, über 64 Jahre!

90%-Konfidenzintervall

Es soll nun ein 90%-Konfidenzintervall für μ bestimmt werden.



Wenn die vorgegebene Wahrscheinlichkeit zwischen zwei Quantile fällt, rückt man auf das äußere Quantil

$$\begin{aligned} & \left[\bar{x} - z(0,95) \cdot \frac{\sigma}{\sqrt{n}}; \bar{x} + z(0,95) \cdot \frac{\sigma}{\sqrt{n}} \right] \\ &= [98 - 2 \cdot 1,65; 98 + 2 \cdot 1,65] \\ &= [98 - 3,3; 98 + 3,3] \\ &= [94,7; 101,3] . \end{aligned}$$

Dieses Intervall ist schmaler als das 95%-Intervall.

7.1.2.

Ausgewählte Konfidenzintervalle

<< hoch zum Anfang vom Inhaltsverzeichnis

<< hoch zum Kapitel 7 vom Inhaltsverzeichnis

Konfidenzintervalle für den Durchschnitt einer Grundgesamtheit Es sei X_1, \dots, X_n eine unabhängige Stichprobe aus der Grundgesamtheit. Der Stichprobenmittelwert ist:

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$$

und die Stichprobenvarianz:

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2 .$$

Die beobachteten Werte dieser Stichprobenfunktionen deuten wir an mit \bar{x} , und s^2 .

Normalverteiltes Merkmal mit bekannter Varianz

Im obigen Beispiel war die Verteilung des Merkmals in der Grundgesamtheit bekannt und normalverteilt und die Varianz σ^2 war bekannt. Man erhält hier das $1-\alpha$ -Konfidenzintervall für μ , den Durchschnitt des Merkmals in der Grundgesamtheit

$$\left[\bar{x} - z\left(1 - \frac{\alpha}{2}\right) \frac{\sigma}{\sqrt{n}}; \bar{x} + z\left(1 - \frac{\alpha}{2}\right) \frac{\sigma}{\sqrt{n}} \right] .$$

Normalverteiltes Merkmal mit unbekannter Varianz

Ist zwar das Merkmal in der Grundgesamtheit normalverteilt, aber die Varianz unbekannt, muss die Varianz des Merkmals durch s^2 geschätzt werden. Damit erhalten wir ein Zufallsintervall das mit Wahrscheinlichkeit $1-\alpha$ den Parameter enthält:

$$P\left(\bar{X} - t\left(1 - \frac{\alpha}{2}; n-1\right) \frac{S}{\sqrt{n}} \leq \mu \leq \bar{X} + t\left(1 - \frac{\alpha}{2}; n-1\right) \frac{S}{\sqrt{n}}\right) = 1 - \alpha .$$

Daraus folgt das $1-\alpha$ -Konfidenzintervall für den Durchschnitt μ des Merkmals in der Grundgesamtheit:

$$\left[\bar{x} - t\left(1 - \frac{\alpha}{2}; n-1\right) \frac{s}{\sqrt{n}}; \bar{x} + t\left(1 - \frac{\alpha}{2}; n-1\right) \frac{s}{\sqrt{n}} \right] .$$

Das Quantil $t\left(1 - \frac{\alpha}{2}; n-1\right)$ kommt jetzt aus einer t-Verteilung mit $n-1$ Freiheitsgraden. Die t-Verteilung hat eine ähnliche Form wie die Normalverteilung, ist aber etwas breiter. In der hier betrachteten Art (zentral) ist sie ebenfalls symmetrisch. Da sie verschiedene Freiheitsgrade hat, ist sie nur für ausgewählte Quantile tabelliert. Es gilt beispielsweise

$$t(0,975;4) = 2,776$$

und

$$t(0,025;4) = -2,776.$$

Merkmal mit unbekannter Verteilung und bekannter Varianz Ist die Verteilung des Merkmals unbekannt, aber die Varianz σ^2 bekannt, kann man für EX des Merkmals X , das Konfidenzintervall

$$\left[\bar{x} - z\left(1 - \frac{\alpha}{2}\right) \frac{\sigma}{\sqrt{n}}; \bar{x} + z\left(1 - \frac{\alpha}{2}\right) \frac{\sigma}{\sqrt{n}} \right].$$

angeben, falls n groß genug ist (Faustregel $n > 30$).

Merkmal mit unbekannter Verteilung und unbekannter Varianz

Sind Verteilung und Varianz des Merkmals unbekannt, kann man für $n > 50$ das Konfidenzintervall für EX angeben als

$$\left[\bar{x} - z\left(1 - \frac{\alpha}{2}\right) \frac{s}{\sqrt{n}}; \bar{x} + z\left(1 - \frac{\alpha}{2}\right) \frac{s}{\sqrt{n}} \right].$$

Konfidenzintervalle für den Anteilswert einer dichotomen Grundgesamtheit

Modell mit Zurücklegen

Die Verteilung eines Merkmals einer dichotomen Grundgesamtheit lässt sich durch das Urnenmodell beschreiben. Man möchte den Anteilswert p , also den Anteil der Kugeln erster Sorte in der Urne bestimmen. Der Anteilswert wird geschätzt durch

$$\hat{p} = \frac{x}{n},$$

worin x der beobachtete Wert der Anzahl X der Kugeln erster Sorte in der Stichprobe ist.

Bei einem Urnenmodell mit Zurücklegen ist X binomialverteilt. Falls n groß genug ist (als Faustregel gilt: $n > 100$ und $n\hat{p}(1 - \hat{p}) \geq 9$), erhält man das $1-\alpha$ -Konfidenzintervall für p durch eine Approximation der Binomialverteilung mit Hilfe der Normalverteilung:

$$\left[\hat{p} - z\left(1 - \frac{\alpha}{2}\right) \sqrt{\frac{\hat{p}(1 - \hat{p})}{n}}; \hat{p} + z\left(1 - \frac{\alpha}{2}\right) \sqrt{\frac{\hat{p}(1 - \hat{p})}{n}} \right].$$

Exakt lässt sich das Konfidenzintervall mit den Verteilungswerten der Binomialverteilung bestimmen. Dafür muß zum Beispiel für eine untere Vertrauensgrenze ein Parameter p_u für die Binomialverteilung bestimmt werden, das so klein ist, daß die Wahrscheinlichkeit aus einer Binomialverteilung mit den Parametern n und p_u gerade x oder mehr Treffer zu erhalten höchstens die eingeräumte Irrtumswahrscheinlichkeit ist.

Für eine Alternative zu diesem Verfahren ist der Zusammenhang der Binomialverteilung mit der Betaverteilung nützlich. Eine untere Vertrauensgrenze für p_u liefert das α -Quantil der Betaverteilung mit den Parametern x und $n - x + 1$. Eine obere Vertrauensgrenze liefert das $1 - \alpha$ -Quantil der Betaverteilung mit den Parametern $x + 1$ und $n - x$. Dabei handelt es sich nicht um zwei verschiedene Methoden, sondern nur um zwei verschiedene Suchverfahren nach einem geeigneten Parameter für die Binomialverteilung, so dass jeweils der einseitige Test für den Parameter der Binomialverteilung nicht zu Ablehnung führt. Weil Quantile der Betaverteilung durch eine Nullstellensuche in der unvollständigen Beta-Funktion bestimmt werden können, ist die Suchstrategie über die Betaverteilung schon dann leicht zugänglich, wenn man einen numerischen Zugang zur unvollständigen Betafunktion und ein allgemeines Verfahren zur Nullstellensuche zu Verfügung hat. Dies kann ein Vorteil gegenüber der Suche nach einem geeigneten Parameter der Binomialverteilung sein, für den das beobachtete x gerade nicht zur Ablehnung führt.

Die exakte Methode über die Suche nach einem geeigneten Parameter der Binomialverteilung so, dass ein einseitiger Test für die Beobachtung x gerade nicht zu Ablehnung führt, ist nur für die Suche nach einer einseitigen Vertrauensgrenze unverfälscht. Ein unverfälschtes zweiseitiges Konfidenzintervall für den Parameter p der Binomialverteilung muss aus einem unverfälschten zweiseitigen Test für den Parameter p abgeleitet werden. Weil die Binomialverteilung außer für $p = 1/2$ nicht symmetrisch ist, genügt es nicht die Irrtumswahrscheinlichkeit α zu gleichen Teilen auf die beiden Enden der Verteilung aufzuteilen.

Modell ohne Zurücklegen

Bei einem Urnenmodell ohne Zurücklegen ist X hypergeometrisch verteilt. Falls die Bedingungen

- $n > \frac{9}{p(1-p)},$
- $n > 100$

- $n/N \leq 0,05$

erfüllt sind, ist die Approximation der hypergeometrischen Verteilung durch die Normalverteilung brauchbar und man erhält das approximative $(1 - \alpha)$ -Konfidenzintervall für θ

$$\left[p - z \left(1 - \frac{\alpha}{2} \right) \sqrt{\frac{p(1-p)}{n}} \sqrt{\frac{N-n}{N-1}} ; p + z \left(1 - \frac{\alpha}{2} \right) \sqrt{\frac{p(1-p)}{n}} \sqrt{\frac{N-n}{N-1}} \right].$$

7.2.

Hypothesentests

[<< hoch zum Anfang vom Inhaltsverzeichnis](#)

[<< hoch zum Kapitel 7 vom Inhaltsverzeichnis](#)

Einführung an Hand eines Beispiels mit Wurstgläsern

Die Firma HappyWurst stellt Wurstwaren her. Sie ist vor allem für ihre delikate Leberwurst in 250g-Gläsern bekannt. Diese werden durch eine Füllanlage mit der noch heißen, flüssigen Masse befüllt. Um Beanstandungen bezüglich der Füllmenge zu vermeiden, füllt man etwas mehr Masse als 250 g ein. Die Füllmenge schwankt immer leicht, aber es wird ein durchschnittliches Füllgewicht von 260g angestrebt. Die Qualitätssicherung soll die Einhaltung dieser Durchschnittsmenge überprüfen.

Überlegung zur Verteilung der Stichprobe

Es ist aber das durchschnittliche Füllgewicht eines Wurstglases unbekannt. Bekannt ist in diesem Beispiel lediglich, daß das Füllgewicht normalverteilt ist mit einer Varianz $\sigma^2 = 64 \text{ [g}^2\text{]}$.

Wie könnte man nun den Durchschnitt ermitteln? Man könnte eine Stichprobe mit z.B. $n = 16$ Beobachtungen ziehen und versuchen, aus dem arithmetischen Mittel \bar{x} auf das durchschnittliche Füllgewicht der Grundgesamtheit zu schließen.

Wir betrachten nun das Füllgewicht eines Wurstglases. Wir bezeichnen es als Zufallsvariable X . Es soll geprüft werden, ob durchschnittlich 260g in einem Glas sind, d.h. ob $EX = 260$ ist.

Beträgt nun tatsächlich der wahre durchschnittliche Absatz der Grundgesamtheit $\mu_0 = 260$, kann man bei einer genügend großen Stichprobe vermuten, daß \bar{x} in der Nähe von μ_0 liegen müßte. Meistens wird \bar{x} in der Nähe von μ_0 liegen, da aber \bar{x} die Realisation einer Zufallsvariablen ist, kann in sehr wenigen Fällen \bar{x} auch extrem weit von μ_0 weg liegen, so daß man dann μ verkehrt einschätzt.

Man könnte aber ein Intervall um μ_0 bestimmen, in dem bei Vorliegen von μ_0 z.B. 95% aller möglichen \bar{x} -Werte liegen, also

$$P(\bar{x}_u \leq \bar{X} \leq \bar{x}_o) = 0,95.$$

Es wird dann eine konkrete Stichprobe genommen. Fällt \bar{x} nicht in dieses Intervall $[\bar{x}_u; \bar{x}_o]$, ist \bar{x} zu weit von μ_0 weg. Man geht dann davon aus, dass $\mu_0 \neq 260$ ist. Damit man dieses Intervall berechnen kann, müssen Informationen über die Verteilung von \bar{x} verfügbar sein.

Ablauf eines Hypothesentests

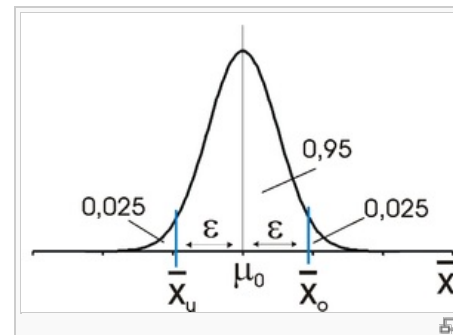
Feststellung der Verteilung des Merkmals in der Grundgesamtheit

Die Zufallsvariable X : Füllgewicht eines Wurstglases ist normalverteilt mit einem unbekannten Erwartungswert μ und der bekannten Varianz $\text{var}X = \sigma^2 = 64$. Man interessiert sich für den Parameter μ .

Aufstellen der Nullhypothese

Man stellt die Nullhypothese $H_0: \mu = \mu_0 = 260$ auf, d.h. man behauptet, das wahre unbekannte durchschnittliche Füllgewicht in der Grundgesamtheit betrage $\mu_0 = 260$.

Festlegen des Nichtablehnungsbereiches für H_0



Zur Überprüfung der Hypothese soll eine Stichprobe im Umfang von $n = 16$ gezogen werden, die zu einer sog. Prüfgröße \bar{X} zusammengefasst wird.

Der Stichprobendurchschnitt \bar{X} ist selbst eine Zufallsvariable und ist als lineare Transformation von X wiederum normalverteilt und zwar mit den Parametern

$$E\bar{X} = \mu \text{ und } var\bar{X} = \frac{\sigma^2}{n}.$$

Bei Gültigkeit von H_0 ist also

$$\bar{X} \rightarrow N\left(\mu_0; \frac{\sigma^2}{n}\right),$$

hier

$$\bar{X} \rightarrow N\left(260; \frac{64}{16} = 4\right).$$

Nun wird der Bereich für \bar{X} festgelegt, in dem die Nullhypothese nicht abgelehnt wird, der Nichtablehnungsbereich (NAB) $[\bar{x}_u; \bar{x}_o]$. Fällt die Prüfgröße \bar{X} in diesem Bereich, wird H_0 nicht abgelehnt. Es soll sein

$$P(\bar{x}_u \leq \bar{X} \leq \bar{x}_o) = 0,95 = 1 - \alpha.$$

Wir nennen α das Signifikanzniveau oder den α -Fehler: Das ist die Wahrscheinlichkeit, dass die Nullhypothese H_0 abgelehnt wird, obwohl $\mu_0 = 260$ der wahre Parameter ist.

Bestimmung von $[\bar{x}_u; \bar{x}_o]$:

Standardisiert man mit

$$Z = \frac{\bar{X} - \mu_0}{\sqrt{\frac{\sigma^2}{n}}},$$

können wir analog zu oben

$$P(z_u \leq Z \leq z_o) = 0,95$$

schreiben. Es ergibt als Intervall für Z :

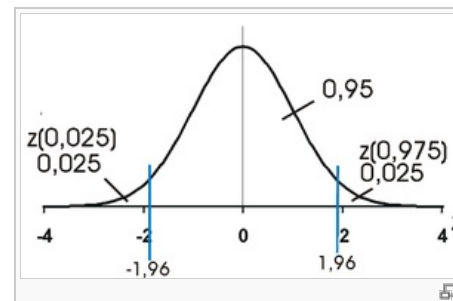
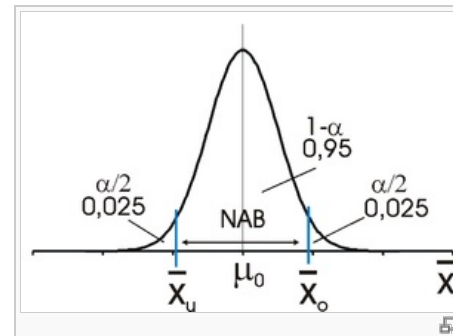
$$\begin{aligned} [z_u; z_o] &= [z(\alpha/2); z(1 - \alpha/2);] \\ &= [-z(1 - \alpha/2); z(1 - \alpha/2);] \\ &= [-z(0,975); z(0,975)] \\ &= [-1,96; 1,96] \end{aligned}$$

Es ist nun aber

$$\bar{x}_u = \mu_0 - z(1 - \alpha/2) \frac{\sigma}{\sqrt{n}} \text{ und } \bar{x}_o = \mu_0 + z(1 - \alpha/2) \frac{\sigma}{\sqrt{n}}$$

so dass hier der Nichtablehnungsbereich für \bar{X}

$$[\bar{x}_u; \bar{x}_o] = [260 - 1,96 \cdot 2; 260 + 1,96 \cdot 2]$$



$$\begin{aligned}
 [\infty; \infty] &= [260 - 3,92; 260 + 3,92] \\
 &= [256,08; 263,92]
 \end{aligned}$$

ist.

Wenn μ_0 tatsächlich 260 ist, würde \bar{x} in 5% aller Stichproben in den Ablehnungsbereich

$$(-\infty; 256,08] \cup [263,92; \infty)$$

fallen.

Stichprobe erheben

Nach der Festlegung des Nichtablehnungsbereichs wird eine Stichprobe genommen. Es wurde hier der Inhalt von 16 Gläsern gewogen. Es ergab sich die Urliste

268 252 254 252 251 245 257 275 268 270

Es ist dann

$$\bar{x} = \frac{1}{16}(268 + 252 + \dots + 267) = \frac{1}{16}(4144) = 259$$

Entscheidung treffen

Wir fällen nun die Entscheidung: Da $\bar{x} = 259$ im Nichtablehnungsbereich liegt, wird H_0 nicht abgelehnt. Es wird davon ausgegangen, dass die Maschine die Gläser korrekt befüllt.

Eine äquivalente Vorgehensweise ist, man bestimmt zunächst die standardisierte Prüfgröße z :

$$z = \frac{\bar{x} - \mu_0}{\frac{\sigma}{\sqrt{n}}} = \frac{259 - 260}{\frac{8}{\sqrt{16}}} = \frac{-1}{2} = -0,5.$$

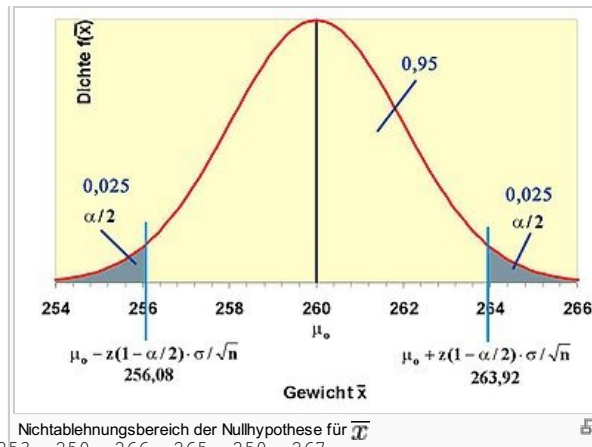
Der Nichtablehnungsbereich für Z ist $[-1,96; 1,96]$. Da z in den Nichtablehnungsbereich fällt, wird H_0 nicht abgelehnt.

Beide Vorgehensweisen liefern das gleiche Ergebnis.

Punkt- und Bereichshypothesen

In obigen Beispiel wurde für das wahre μ nur ein bestimmter Punkt getestet: $H_0: \mu = \mu_0$, also handelt es sich um eine Punkthypothese. Es könnte aber sein, dass der Hersteller einem Großabnehmer versichert hat, dass das durchschnittliche Füllgewicht mindestens 260 g beträgt. Es wird also hier genügen, zu prüfen, ob der Mindestwert erreicht wird. Es ist aber kein Problem, wenn die durchschnittliche Füllmenge größer als 260 ist.

Man stellt also als Arbeitshypothese auf $H_1: \mu > \mu_0 = 260$. Wird die



man stellt also die Alternativhypothese auf: $H_0: \mu \geq \mu_0 = 260$. Wird die Prüfgröße \bar{x} geringfügig kleiner als 260, kann das eine Zufallsschwankung sein. Aber wird \bar{x} zu klein, muss H_0 abgelehnt werden. Da hier nur der Bereich links von μ_0 kritisch für die Ablehnung ist, wird das gesamte α links auf dem Zahlenstrahl platziert, der kritische Wert für z ist also $z(\alpha) = -z(1-\alpha)$. Fällt z in den Ablehnungsbereich $(-\infty; -z(1-\alpha)]$, wird H_0 abgelehnt. Man geht dann davon aus, dass μ kleiner als μ_0 sein muss, dass also die Befüllung nicht ordnungsgemäß ist. Der kritische Wert für \bar{x} ist hier

$$\bar{x}_{1-\alpha} = \mu_0 - z(1-\alpha) \cdot \frac{\sigma}{\sqrt{n}},$$

also

$$\bar{x}_{1-\alpha} = 260 - 1,65 \cdot \frac{8}{\sqrt{16}} = 256,7.$$

Wenn die Stichprobe ein Durchschnittsgewicht von weniger als 256,7g ergibt, wird die Lieferung beanstandet.

Entsprechend erhält man unter der Hypothese $H_0: \mu \leq \mu_0$ für die Prüfgröße z den Ablehnungsbereich $[z(1-\alpha); \infty)$ bzw.

$$\bar{x}_{1-\alpha} = \mu_0 + z(1-\alpha) \cdot \frac{\sigma}{\sqrt{n}}.$$

Fehler und Varianzen

Fehlerarten

Warum wird der α -Fehler als Fehler bezeichnet? Hier wollen wir uns zunächst mal überlegen, welche Fehler bei der Entscheidung überhaupt gemacht werden können?

1. H_0 ist wahr, die Prüfgröße fällt aber in den Ablehnungsbereich (in $\alpha \cdot 100\%$ aller Stichproben). Hier würde man H_0 irrtümlicherweise ablehnen, obwohl H_0 wahr ist: α -Fehler oder Fehler 1. Art. In unserem Beispiel würde also die Lieferung möglicherweise zurückgewiesen werden, obwohl die Gläser korrekt befüllt worden sind.

1. H_0 ist falsch, die Prüfgröße fällt aber in den Nichtablehnungsbereich. In Wirklichkeit ist $\mu = \mu_1$, z.B. $\mu_1 = 255$ g. Jetzt ist bei unveränderter Varianz in Wahrheit der Stichprobendurchschnitt \bar{x} verteilt wie

$$N(\mu_1; \frac{\sigma^2}{n}) = N(255; 4)$$

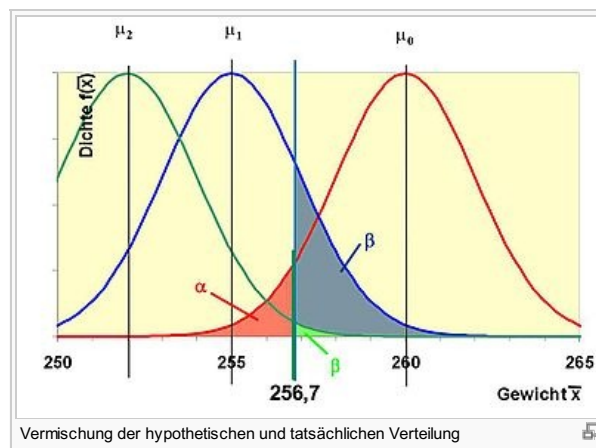
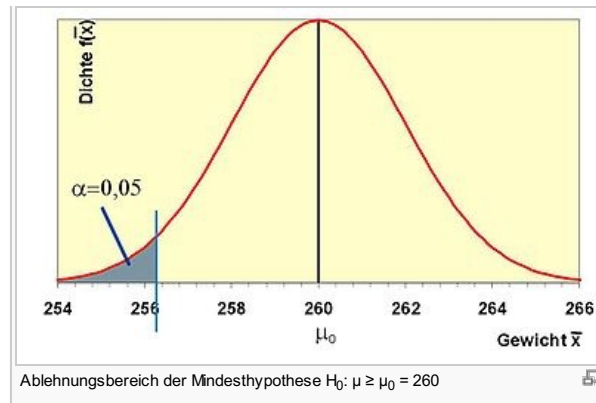
Unter dieser Verteilung beträgt die Wahrscheinlichkeit, dass H_0 (fälschlicherweise) nicht abgelehnt wird,

$$P(\bar{X} \geq 256,7) = 1 - \Phi_{\bar{x}}(256,7|255; 4),$$

was sich einfach berechnen lässt als

$$1 - \Phi_z\left(\frac{256,7 - 255}{2}\right) = 1 - \Phi_z(0,85) = 0,1977.$$

Man würde also mit fast 20%iger Wahrscheinlichkeit irrtümlicherweise die Lieferung akzeptieren. Dieser Fehler ist der β -Fehler oder Fehler 2.



Art.

Wenn in Wahrheit $\mu = \mu_2 = 252$ ist, beträgt der β -Fehler

$$P(\bar{X} \geq 256,7) = 1 - \Phi_{\bar{x}}(256,7|252;4) = 1 - \Phi_z\left(\frac{256,7 - 252}{2}\right) = 1 - \Phi_z(2,35) = 0,0094.$$

Hier ist die Wahrscheinlichkeit einer irrtümlichen Ablehnung schon sehr klein.

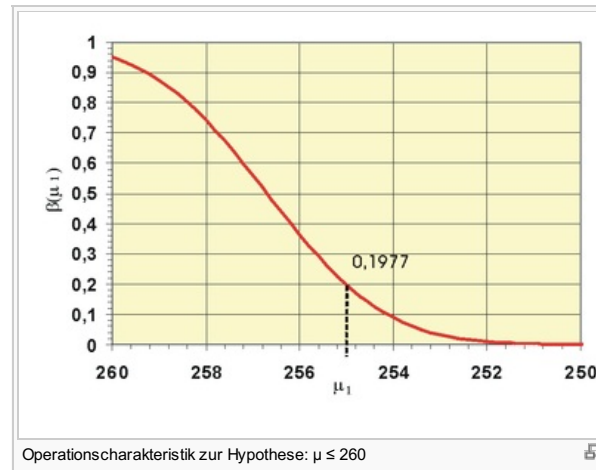
Der β -Fehler hängt also von μ_1 ab. Man kann den β -Fehler in

Abhängigkeit von μ_1 als Funktion darstellen: $\beta = f(\mu_1)$. Diese Funktion nennt man Operationscharakteristik. Der Wert $1 - \beta$ ist dagegen die Wahrscheinlichkeit, dass H_0 abgelehnt wird, wenn μ_1 der wahre Parameter ist. Man sieht an der Grafik, dass $1 - \beta$ für $\mu = 260$ gerade 0,05 ist. Das ist natürlich die Wahrscheinlichkeit, dass H_0 (hier fälschlicherweise) abgelehnt wird, wenn 260 tatsächlich der wahre Parameter ist.

Um die Wahrscheinlichkeit für eine falsche Entscheidung zu reduzieren, ist es wünschenswert, möglichst schnell in den Bereich $\beta \approx 0$ zu kommen. U. U. hilft eine Erhöhung des Stichprobenumfangs.

Eine Hypothese, die nicht abgelehnt ist, gilt nicht automatisch als angenommen, denn der β -Fehler ist i.a. unbekannt.

Wenn ein Test die Wahrscheinlichkeit der Annahme falscher Nullhypothesen möglichst reduziert, nennt man ihn trennscharf.



Breite des Nichtablehnungsbereichs

Es soll nun wieder die Punkthypothese $H_0: \mu = \mu_0$ betrachtet werden. Es ergab sich hier für \bar{x} der Nichtablehnungsbereich $[256,08; 263,92]$ mit einer Breite 7,84 g.

Änderung des Signifikanzniveaus

Welcher NAB ergibt sich für $\alpha = 0,01$? Wir errechnen das $(1 - \alpha/2)$ -Quantil als

$$\alpha = 0,01 \rightarrow \alpha/2 = 0,005 \rightarrow 1 - \alpha/2 = 0,995 \rightarrow z(0,995) = 2,58$$

und erhalten den Nichtablehnungsbereich für \bar{x} als

$$\begin{aligned} & [260 - z(0,995) \cdot 2; 260 + z(0,995) \cdot 2] \\ &= [260 - 2,58 \cdot 2; 260 + 2,58 \cdot 2] \\ &= [260 - 5,16; 260 + 5,16] \\ &= [254,84; 265,16] \end{aligned}$$

Hier ist der Nichtablehnungsbereich breiter als für $\alpha = 0,05$: H_0 wird nur in 1% aller Stichproben fälschlicherweise abgelehnt. Hier hätte die Lieferfirma einen Vorteil.

Welcher NAB ergibt sich für $\alpha = 0,1$?

$$\begin{aligned} & [260 - z(0,95) \cdot 2; 260 + z(0,95) \cdot 2] \\ &= [260 - 1,65 \cdot 2; 260 + 1,65 \cdot 2] \\ &= [260 - 3,30; 260 + 3,30] \\ &= [256,70; 263,30] \end{aligned}$$

Hier ist der Nichtablehnungsbereich schmaler, H_0 wird in 10% aller Stichproben fälschlicherweise abgelehnt.

Änderung der Varianz

Was passiert, wenn die Varianz $\sigma^2 = 256$ ist ($\alpha = 0,05$)? Man erhält hier für die Punkthypothese $H_0: \mu = \mu_0 = 260$ den NAB für \bar{x}

$$\begin{aligned} & [260 - 1,96 \cdot \sqrt{\frac{256}{16}}; 260 + 1,96 \cdot \sqrt{\frac{256}{16}}] \\ &= [260 - 1,96 \cdot 4; 260 + 1,96 \cdot 4] \\ &= [260 - 7,84; 260 + 7,84] \\ &= [252,16; 267,84] \end{aligned}$$

Die Breite des Nichtablehnungsbereichs ist hier 15,68g.

Für $H_1: \mu \geq \mu_0$ ergibt sich dann entsprechend als kritischer Wert

$$260 - 1,65 \cdot 4 = 253,4$$

Die Grafik zeigt den Fall der Bereichshypothese mit einer Varianz von 16: Durch die große Varianz sind die Normalverteilungskurven sehr flach und durchmischen sich stark. Der Betafehler bei $\mu_1 = 255$ ist sehr groß. Eine vernünftige Kontrolle der Abfüllmaschine ist nicht mehr möglich.

Der Nichtablehnungsbereich wird mit wachsender Varianz breiter, der Test verliert an Trennschärfe.

Änderung des Stichprobenumfangs

Was passiert, wenn der Stichprobenumfang jetzt 64 beträgt ($\alpha = 0,05$; $\sigma^2 = 64$)?

$$\begin{aligned} & [260 - 1,96 \cdot \sqrt{\frac{64}{64}}; 260 + 1,96 \cdot \sqrt{\frac{64}{64}}] \\ &= [260 - 1,96 \cdot 1; 260 + 1,96 \cdot 1] \\ &= [260 - 1,96; 260 + 1,96] \\ &= [258,04; 261,96] \end{aligned}$$

Hier hat der Nichtablehnungsbereich eine Breite von 3,92, denn durch den größeren Stichprobenumfang hat sich die Varianz von \bar{x} verringert. Der NAB schrumpft bei steigendem Stichprobenumfang, der Test wird trennschärfer.

Mindest erforderlicher Stichprobenumfang

Wie groß muß die Stichprobe mindestens sein, damit die Breite des NAB für $\alpha = 0,05$ höchstens 10 beträgt?

Die Breite des NAB ist ja definiert durch

$$2 \cdot z(1 - \alpha/2) \cdot \frac{\sigma}{\sqrt{n}}$$

Es soll also hier sein

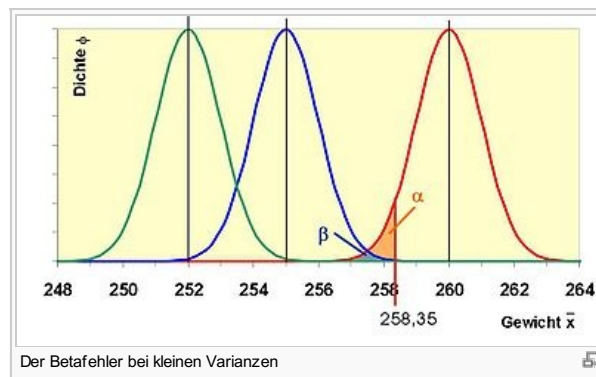
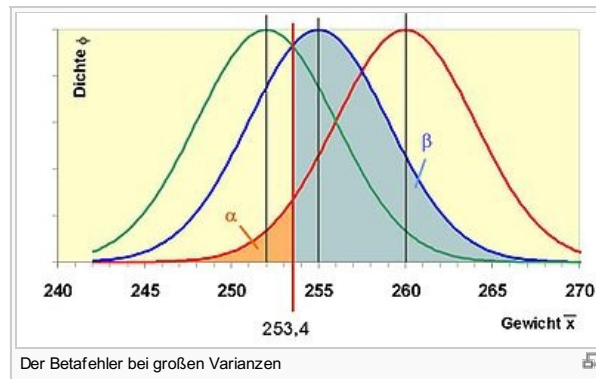
$$2 \cdot 1,96 \cdot \frac{64}{\sqrt{n}} \leq 10$$

Die Auflösung der Ungleichung nach \sqrt{n} ergibt dann

$$\sqrt{n} \geq 2 \cdot 1,96 \cdot \frac{64}{10} = 25,088 \text{ und } (\sqrt{n})^2 = 629,41$$

Da wir nur ganze Wurstgläser analysieren können, brauchen wir einen Stichprobenumfang von mindestens 630 Gläsern.

Kann die Wurst mit dem Glas zusammen gewogen werden, stellt diese hohe Zahl kein Problem dar. Geht durch so eine Stichprobe allerdings die Zerstörung der Ware mit einher, etwa die lebensmitteltechnische Untersuchung einer Konservendose, muss man einen Kompromiss



zwischen mangelnder Trennschärfe und Zerstörung der Ware finden.

7.2.1.	Verteilungsgebundene Tests	<< hoch zum Anfang vom Inhaltsverzeichnis << hoch zum Kapitel 7 vom Inhaltsverzeichnis
7.2.1.1.	Test auf Erwartungswert	<< hoch zum Anfang vom Inhaltsverzeichnis << hoch zum Kapitel 7 vom Inhaltsverzeichnis

Erwartungswert

1. Bekannte Verteilung und Varianz

Im einführenden Beispiel war die Art der Verteilung des Merkmals in der Grundgesamtheit bekannt, namentlich eine Normalverteilung mit bekannter Varianz. Die Prüfgröße

$$z = \frac{\bar{x} - \mu_0}{\frac{\sigma}{\sqrt{n}}}$$

ist dann unter der Nullhypothese $\mu = \mu_0$ standardnormalverteilt. Wir erhalten die Entscheidungsregeln für eine gewählte Irrtumswahrscheinlichkeit α

- $H_0: \mu = \mu_0$ wird abgelehnt, falls $z < -z(1-\alpha/2)$ oder $z > z(1-\alpha/2)$ ist.
- $H_0: \mu \leq \mu_0$ wird abgelehnt, falls $z > z(1-\alpha)$ ist.
- $H_0: \mu \geq \mu_0$ wird abgelehnt, falls $z < -z(1-\alpha)$ ist.

2. Bekannte Verteilung und unbekannte Varianz

Häufig wird neben dem Erwartungswert die Varianz ebenfalls nicht bekannt sein, so dass man statt der Varianz in der Grundgesamtheit die Schätzung

$$\hat{\sigma}^2 = s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

verwendet. Wir erhalten nun bei normalverteilter Grundgesamtheit statt z die Prüfgröße

$$t = \frac{\bar{x} - \mu_0}{\frac{s}{\sqrt{n}}},$$

die t-verteilt mit $n-1$ Freiheitsgraden ist.

Die t-Verteilung hat eine ähnliche Form wie die Normalverteilung. In der hier betrachteten Art (zentrale t-Verteilung) ist sie ebenfalls symmetrisch bezüglich der Null. Da sie verschiedene Freiheitsgrade hat, ist sie nur für ausgewählte Quantile tabelliert. Es ist $t(p;k)$ das p-Quantil der t-Verteilung mit k Freiheitsgraden.

Es gilt beispielsweise für die Zufallsvariable t mit 5 Freiheitsgraden:

$$P(t \leq 3,365) = 0,99 \text{ bzw. } t(0,99;5) = 3,365.$$

Wir erhalten die Entscheidungsregeln

- $H: \mu = \mu_0$ wird abgelehnt, falls $t < -t(1-\alpha/2; n-1)$ oder $t > t(1-\alpha/2; n-1)$ ist.
- $H: \mu \leq \mu_0$ wird abgelehnt, falls $t > t(1-\alpha; n-1)$ ist.
- $H: \mu \geq \mu_0$ wird abgelehnt, falls $t < -t(1-\alpha; n-1)$ ist.
- Ist $n > 30$, können die Quantile der t-Verteilung durch die entsprechenden Quantile der Normalverteilung ersetzt werden.

3. Unbekannte Verteilung und bekannte Varianz

Ist die Verteilung des Merkmals X unbekannt, aber die Varianz $\text{var}X$ bekannt, verwendet man bei einem $n > 30$ die standardnormalverteilte

Prüfgröße

$$z = \frac{\bar{x} - \mu_0}{\frac{s}{\sqrt{n}}}$$

Wir erhalten die Entscheidungsregeln analog zu 1.

4. Unbekannte Verteilung und unbekannte Varianz

Sind Verteilung und Varianz des Merkmals X in der Grundgesamtheit unbekannt, verwendet man für $n > 50$ die standardnormalverteilte Prüfgröße

$$z = \frac{\bar{x} - \mu_0}{\frac{s}{\sqrt{n}}}$$

Wir verwenden die Entscheidungsregeln analog zu 1.

7.2.1.2.

Test auf Anteilswert

[<< hoch zum Anfang vom Inhaltsverzeichnis](#)

[<< hoch zum Kapitel 7 vom Inhaltsverzeichnis](#)

Anteilswert einer dichotomen Grundgesamtheit

Die Verteilung des Merkmals X einer dichotomen Grundgesamtheit lässt sich durch das Urnenmodell beschreiben. Man möchte den Anteilswert θ , also den Anteil der Kugeln erster Sorte in der Urne bestimmen. Der Anteilswert wird geschätzt durch

$$\hat{\theta} = p = \frac{x}{n},$$

wobei x die Zahl der Kugeln erster Sorte in der Stichprobe ist. Bei einem Urnenmodell mit Zurücklegen ist X binomialverteilt.

Falls

$$n > \frac{9}{\theta \cdot (1 - \theta)}$$

können wir die Prüfgröße verwenden

$$z = \frac{x \pm 0,5 - n \cdot \theta_0}{\sqrt{n \cdot \theta \cdot (1 - \theta)}}$$

- $H_0: \theta = \theta_0$ wird abgelehnt, falls

$$z = \frac{x + 0,5 - n \cdot \theta_0}{\sqrt{n \cdot \theta \cdot (1 - \theta)}} < -z(1 - \alpha/2),$$

(wenn die Prüfgröße $z < 0$ ist) oder

$$z = \frac{x - 0,5 - n \cdot \theta_0}{\sqrt{n \cdot \theta \cdot (1 - \theta)}} > z(1 - \alpha/2)$$

(wenn die Prüfgröße $z > 0$ ist) errechnet wird.

- $H_0: \theta \leq \theta_0$ wird abgelehnt, falls

$$z > z = \frac{x - 0,5 - n \cdot \theta_0}{\sqrt{n \cdot \theta \cdot (1 - \theta)}} > z(1 - \alpha)$$

ist.

- $H_0: \theta \geq \theta_0$ wird abgelehnt, falls

$$z = \frac{x + 0,5 - n \cdot \theta_0}{\sqrt{n \cdot \theta \cdot (1 - \theta)}} < -z(1 - \alpha)$$

ist.

Ist n zu klein, kann der Ablehnungsbereich mit Hilfe der F-Verteilung exakt bestimmt werden oder mit dem Prinzip des konservativen Testens festgelegt werden.

7.2.1.3.

Test auf Varianz

[<< hoch zum Anfang vom Inhaltsverzeichnis](#)

[<< hoch zum Kapitel 7 vom Inhaltsverzeichnis](#)

Test auf Varianz

Herleitung der Prüfgröße Betrachten wir eine normalverteilte Grundgesamtheit. Die Schätzung für die Varianz ist hier

$$\hat{\sigma}^2 = s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2.$$

Wir wollen nun eine passende Prüfgröße für einen Varianztest herleiten. Seien $x_i : i = 1, \dots, n$ unabhängige, normalverteilte Zufallsgrößen mit Mittelwert μ und Varianz σ^2 . Dann sind die Größen $(x_i - \mu) / \sigma$ unabhängige, standardnormalverteilte Zufallsvariablen und die Summe der Quadrate ist χ^2 -verteilt mit n Freiheitsgraden:

$$\sum_{i=1}^n \frac{(x_i - \mu)^2}{\sigma^2}.$$

Schätzt man

$$\hat{\mu} = \bar{x}$$

geht ein Freiheitsgrad verloren.

$$\sum_{i=1}^n \frac{(x_i - \bar{x})^2}{\sigma^2}$$

ist χ^2 -verteilt mit $n-1$ Freiheitsgraden. Wir wollen nun diese Summe mit S^2 verquicken, um eine Prüfgröße für diesen Test zu erhalten. Es ist dann

$$\sum_{i=1}^n \frac{(x_i - \bar{x})^2}{n-1} \cdot (n-1) = \frac{S^2 \cdot (n-1)}{\sigma^2}$$

ebenfalls χ^2 -verteilt mit $n-1$ Freiheitsgraden. Unter der Nullhypothese $H_0: \sigma^2 = \sigma_0^2$ ist dann

$$Y = \frac{S^2 \cdot (n-1)}{\sigma_0^2}$$

ebenfalls verteilt wie oben.

Wir wollen nun für $H_0: \sigma^2 = \sigma_0^2$ den Nichtablehnungsbereich für den Test angeben. Die Hypothese wird nicht abgelehnt, wenn die Prüfgröße y in das Intervall

$$[\chi^2(\frac{\alpha}{2}; n-1); \chi^2(1 - \frac{\alpha}{2}; n-1)]$$

fällt, wobei $\chi^2(p;k)$ das p -Quantil der χ^2 -Verteilung mit k Freiheitsgraden ist.

Die Nichtablehnungsbereiche für die Bereichshypothesen werden analog zu der Vorgehensweise bei Erwartungswerten festgelegt:

Bei der Mindesthypothese $H_0 : \sigma^2 \geq \sigma_0^2$ wird die Hypothese abgelehnt, wenn die Prüfgröße

$$Y < \chi^2(\alpha; n-1) \text{ ist.}$$

Bei der Höchsthypothese $H_0 : \sigma^2 \leq \sigma_0^2$ wird die Hypothese abgelehnt, wenn die Prüfgröße

$$Y > \chi^2(1-\alpha; n-1) \text{ ist.}$$

Beispiel für eine Punkthypothese

Ein großer Blumenzwiebelzüchter hat eine neue Sorte von Lilien gezüchtet. Die Zwiebeln sollen im Verkauf in verschiedenen Größenklassen angeboten werden. Um das Angebot planen zu können, benötigt der Züchter eine Information über die Varianz der Zwiebelgröße. Es wurden 25 Zwiebeln zufällig ausgewählt und gemessen. Man erhielt die Durchmesser (cm)

8 10 9 7 6 10 8 8 8 6 7 9 7 10 9 6 7 7 8 8 8 10 10 7 7

Es soll die Hypothese überprüft werden, dass die Varianz der Zwiebelgröße 3 cm² beträgt ($\alpha = 0,05$).

Die Nullhypothese lautet $H_0 : \sigma^2 = \sigma_0^2 = 3$

Nichtablehnungsbereich für die Prüfgröße y ist

$$[\chi^2(\frac{\alpha}{2}; n-1); \chi^2(1-\frac{\alpha}{2}; n-1)] =$$

$$[\chi^2(0,025; 24); \chi^2(0,975; 24)] = [12,40; 39,36].$$

Es ergab sich für die Stichprobe $\bar{x} = 8$ und $s^2 = \frac{42}{24} = 1,75$. Die Prüfgröße errechnet sich als

$$y = \frac{S^2 \cdot (n-1)}{\sigma_0^2} = \frac{1,75 \cdot 24}{3} = 14.$$

Die Hypothese kann nicht abgelehnt werden.

Beispiel für eine Bereichshypothese

An einer Abfüllanlage werden Tagesdosen für ein sehr teures flüssiges Medikament in Plastischälchen eingebracht. Da das Medikament hochwirksam ist, soll die Abweichung der Füllmenge vom Mittelwert möglichst wenig schwanken. Man weiß, dass die Füllmenge normalverteilt ist. Zur Kontrolle soll die Hypothese getestet werden, dass die Varianz höchstens 0,01 ml² beträgt. Eine Stichprobe von 20 Schälchen ergab den Mittelwert 0,5 und die Varianz 0,014.

Zu testen ist $H_0 : \sigma^2 \leq \sigma_0^2$.

$$\text{Die Prüfgröße für } H_0 \text{ ist } Y = \frac{S^2 \cdot (n-1)}{\sigma_0^2}.$$

Die Hypothese wird abgelehnt, wenn $y > \chi^2(1-\alpha; n-1) = \chi^2(0,9; 19) = 27,20$ ist.

Die Stichprobe ergab

$$y = \frac{0,014 \cdot 19}{0,01} = 26,6$$

Die Hypothese wird nicht abgelehnt. Man geht davon aus, dass die Varianz der Füllmenge sich nicht verändert hat.

Vergleich zweier Varianzen

Wir haben es mit zwei verschiedenen Grundgesamtheiten zu tun. Wir interessieren uns dafür, ob die Varianzen dieser beiden Grundgesamtheiten gleich sind. Beide Merkmale dieser Grundgesamtheiten sollen normalverteilt sein.

Herleitung der Prüfgröße

Zu prüfen ist also die Hypothese: $H_0: \sigma_1^2 = \sigma_2^2$.

Geschätzt werden beide Varianzen wieder mit der Stichprobenvarianz

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2.$$

Es soll nun daraus eine Prüfgröße konstruiert werden. Wir wissen bereits, dass der Quotient

$$Y = \sum_{i=1}^n \frac{(X_i - \bar{X})^2}{\sigma^2}$$

χ^2 -verteilt mit $n-1$ Freiheitsgraden ist. Eine Möglichkeit, zwei solche Zufallsvariablen zu verquicken, ist die F-Verteilung. Es ist nämlich der Quotient

$$f = \frac{\frac{Y_1}{n_1-1}}{\frac{Y_2}{n_2-1}} = \frac{\frac{\sum_{i=1}^{n_1} (X_{1i} - \bar{X}_1)^2}{(n_1-1)\sigma_1^2}}{\frac{\sum_{i=1}^{n_2} (X_{2i} - \bar{X}_2)^2}{(n_2-1)\sigma_2^2}}$$

F-verteilt mit $n_1 - 1$ und $n_2 - 1$ Freiheitsgraden. Wir müssen nun noch unsere Stichprobenvarianzen einpflegen und wir sehen, dass ja in Zähler und Nenner die Stichprobenvarianzen S_1^2 und S_2^2 schon dastehen. Also erhalten wir

$$f = \frac{\frac{S_1^2}{\sigma_1^2}}{\frac{S_2^2}{\sigma_2^2}} = \frac{S_1^2}{S_2^2} \cdot \frac{\sigma_2^2}{\sigma_1^2}$$

Wir wollen diesen Quotienten nun mit der Nullhypothese in Verbindung bringen. Die Hypothese

$H_0: \sigma_1^2 = \sigma_2^2$ lässt sich auch schreiben als $H_0: \frac{\sigma_1^2}{\sigma_2^2} = 1$ und es ist dann der Quotient der Prüfgröße unter H_0

$$f = \frac{S_1^2}{S_2^2} \cdot 1.$$

Wenn die Nullhypothese wahr ist, sollte f nicht zu groß sein, aber auch nicht zu klein, weil sonst die Stichprobenvarianzen zu unterschiedlich wären. H_0 wird also nicht abgelehnt, wenn die Stichprobe f in den „mittleren“ Bereich

$$\left[f\left(\frac{\alpha}{2}; n_1 - 1; n_2 - 1\right); f\left(1 - \frac{\alpha}{2}; n_1 - 1; n_2 - 1\right) \right]$$

fällt, wobei $f(p; k_1; k_2)$ das p -Quantil der F-Verteilung mit k_1 und k_2 Freiheitsgraden ist.

Bereichshypothesen werden entsprechend aufgefasst:

$H_0: \sigma_1^2 \leq \sigma_2^2$ lässt sich auch schreiben als $H_0: \frac{\sigma_1^2}{\sigma_2^2} \leq 1$.

Dieser Test wird abgelehnt, wenn

$$f > f\left(1 - \frac{\alpha}{2}; n_1 - 1; n_2 - 1\right),$$

wobei sich f wie oben berechnet.

Entsprechend wird $H_0 : \frac{\sigma_1^2}{\sigma_2^2} \geq 1$ abgelehnt, wenn

$$f < f\left(\frac{\alpha}{2}; n_1 - 1; n_2 - 1\right).$$

Beispiel

Bert und Berta haben im Fach Analysis ein Tutorium gehalten. Die Zeit, die die n_1 bzw. n_2 Studierenden für eine typische Klausuraufgabe benötigten, wurde festgehalten:

Tutorium von Bert: 8 3 4 4 10 9 2 9
Tutorium von Berta: 5 4 7 6 4

Beide Gruppen erzielten eine durchschnittliche Bearbeitungsdauer von 6 min. Ist aber auch die Varianz beider Gruppenleistungen gleich?

Wir wollen also nun bei einem Signifikanzniveau 0,05 die Nullhypothese testen, dass die Varianzen gleich sind.

Der Nichtablehnungsbereich für diesen Test ist

$$\begin{aligned} & [f(\frac{\alpha}{2}; n_1 - 1; n_2 - 1); f(1 - \frac{\alpha}{2}; n_1 - 1; n_2 - 1)] \\ &= [f(0,025; 8; 5); f(0,975; 8; 5)] \\ &= [0,21; 6,76] \end{aligned}$$

wobei sich

$$f(0,025; 8; 5) = \frac{1}{f(0,975; 5; 8)} = \frac{1}{4,82} = 0,21$$

errechnet. Wir erhalten zunächst die Stichprobenvarianzen

$$\begin{aligned} s_1^2 &= \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 \\ &= \frac{1}{8} ((8-6)^2 + (3-6)^2 + (4-6)^2 + \dots + (9-6)^2) = \frac{72}{8} = 9 \end{aligned}$$

und analog dazu

$$s_2^2 = 5,2.$$

Die Prüfgröße errechnet sich nun als

$$\frac{S_1^2}{S_2^2} \cdot 1 = \frac{9}{5,2} \cdot 1 = 1,73$$

Sie fällt in den Nichtablehnungsbereich und man kann die Hypothese nicht ablehnen.

7.2.2.	Verteilungsfreie Tests	<< hoch zum Anfang vom Inhaltsverzeichnis << hoch zum Kapitel 7 vom Inhaltsverzeichnis
7.2.2.1.	Prüfung des Zusammenhangs zweier Merkmale	<< hoch zum Anfang vom Inhaltsverzeichnis << hoch zum Kapitel 7 vom Inhaltsverzeichnis

Stochastische Unabhängigkeit

Die Beobachtungen zweier Merkmale X und Y liegen als gemeinsame klassierte Häufigkeitsverteilung vor mit n und m Kategorien und den dazugehörigen gemeinsamen Häufigkeiten n_{ij} ($i = 1, \dots, n; j = 1, \dots, m$) vor. Zur Prüfung der Hypothese H_0 : „X und Y sind stochastisch unabhängig“ verwendet man die Prüfgröße

$$\chi^2 = \sum_i \sum_j \frac{(n_{ij} - \frac{n_{i.} n_{.j}}{n})^2}{\frac{n_{i.} n_{.j}}{n}}$$

Es soll jedes $\frac{n_{i.} n_{.j}}{n} \geq 5$ sein. Falls diese Forderung nicht gegeben ist, müssen so viele Zeilen und/oder Spalten zusammengefasst werden, bis die Vorgabe erfüllt ist.

Die Hypothese, dass X und Y stochastisch unabhängig sind, wird abgelehnt, wenn $\chi^2 > \chi^2(1 - \alpha; (m - 1)(n - 1))$ ist, als $(1 - \alpha)$ -Quantil der χ^2 -Verteilung mit $(m - 1)(n - 1)$ Freiheitsgraden.

Bemerkung: Dieser Test kann auch für die Prüfung der stochastischen Unabhängigkeit zweier Ereignisse verwendet werden. Man spricht hier von einem Vierfelder-Test.

Korrelation

Normalverteilung beider Merkmale

Die Merkmale X und Y sind normalverteilt. Es wird die spezielle Nullhypothese $H_0: \rho_{xy} = 0$ geprüft. Man schätzt den Korrelationskoeffizienten ρ mit dem Korrelationskoeffizienten r nach Bravais-Pearson und verwendet die Prüfgröße

$$t = \frac{r}{\sqrt{\frac{1-r^2}{n-2}}}$$

$H_0: \rho_{xy} = 0$ wird abgelehnt, falls $t < -t(1-\alpha/2; n - 2)$ oder $t > t(1-\alpha/2; n - 2)$ ist.

Wird H_0 abgelehnt, geht man davon aus, dass X und Y korreliert sind. Sie sind dann auch stochastisch abhängig, so dass dieser Test im Ablehnungsfall auch die stochastische Unabhängigkeit erfasst. Bei Nichtablehnung können die Merkmale trotzdem abhängig sein, denn der Korrelationskoeffizient misst bekanntlich nur die lineare Abhängigkeit.

Wird $H_0: \rho_{xy} = \rho_0 \neq 0$ geprüft, hat r eine sog. nichtzentrale Verteilung, die nicht mehr ohne weiteres berechnet werden kann und nur noch näherungsweise mit der sog. Fisherschen Transformation angebar ist.

Unbekannte Verteilung beider Merkmale

Die Merkmale X und Y sind beliebig verteilt. Es wird die spezielle Nullhypothese $H_0: \rho_{xy} = 0$ geprüft. Man schätzt den Korrelationskoeffizienten ρ mit dem Rangkorrelationskoeffizienten nach Spearman-Pearson r_{SP} .

Für $n > 10$ verwendet man die Prüfgröße

$$t = \frac{r_{SP}}{\sqrt{\frac{1-r_{SP}^2}{n-2}}}$$

$H_0: \rho_{xy} = 0$ wird abgelehnt, falls $t < -t(1-\alpha/2; n - 2)$ oder $t > t(1-\alpha/2; n - 2)$ ist.

Parameter der linearen Regression

Ausgegangen wird von der unbekannten Regressionsgeraden

$$y = \alpha + \beta x + u$$

und der Schätzung

$$y = a + bx + d.$$

Die Störgröße u ist normalverteilt:

$$u \rightarrow N(0; \sigma^2).$$

Die Varianz der Störgröße σ^2 wird geschätzt mit

$$s^2 = \frac{1}{n-2} \sum_i (d_i - \bar{d})^2 = \frac{1}{n-2} \sum_i d_i^2$$

Es ist auch

$$\sum_i d_i^2 = (1 - r^2) \cdot \sum_i (y_i - \bar{y})^2$$

Steigungskoeffizient β

β wird geschätzt durch b . Unter H_0 ist $b \rightarrow N(\beta; \frac{\sigma^2}{\sum_i (x_i - \bar{x})^2})$.

Verwendet wird die Prüfgröße

$$t = \frac{b - \beta_0}{\frac{s}{\sqrt{\sum_i (x_i - \bar{x})^2}}}$$

die unter H_0 t-verteilt ist mit $n-2$ Freiheitsgraden.

- $H_0: \beta = \beta_0$ wird abgelehnt, falls $t < -t(1-\alpha/2; n-2)$ oder $t > t(1-\alpha/2; n-2)$ ist.
- $H_0: \beta \leq \beta_0$ wird abgelehnt, falls $t > t(1-\alpha/2; n-2)$ ist.
- $H_0: \beta \geq \beta_0$ wird abgelehnt, falls $t < -t(1-\alpha/2; n-2)$ ist.

In der Praxis wird meistens $H_0: \beta = 0$ getestet. Wird die Hypothese nicht abgelehnt, scheint x unerheblich für die Erklärung von y zu sein.

Absolutglied α

α wird geschätzt durch a . Unter H_0 ist

$$a \rightarrow N(\alpha_0; \frac{\sigma^2 \cdot \sum_i x_i^2}{\sum_i (x_i - \bar{x})^2})$$

Für den Test verwendet man die Prüfgröße

$$t = \frac{a - \alpha_0}{s} \sqrt{\frac{\sum_i (x_i - \bar{x})^2}{\sum_i x_i^2}},$$

die unter H_0 t-verteilt ist mit $n-2$ Freiheitsgraden.

- $H_0: \alpha = \alpha_0$ wird abgelehnt, falls $t < -t(1-\alpha/2; n-2)$ oder $t > t(1-\alpha/2; n-2)$ ist.
- $H_0: \alpha \leq \alpha_0$ wird abgelehnt, falls $t > t(1-\alpha/2; n-2)$ ist.
- $H_0: \alpha \geq \alpha_0$ wird abgelehnt, falls $t < -t(1-\alpha/2; n-2)$ ist.

In der Praxis wird meistens $H_0: \alpha = 0$ getestet. Wird die Hypothese nicht abgelehnt, geht die wahre Regressionsgerade möglicherweise durch den Nullpunkt des Koordinatensystems.

8.

Übungsaufgaben

[<< hoch zum Anfang vom Inhaltsverzeichnis](#)

[<< hoch zum Kapitel 8 vom Inhaltsverzeichnis](#)

Kapitel 1

Aufgabe 1.1 - Mischanlage für Porzellanfabrik

Eine Porzellanfabrik erhält eine neue Mischanlage für spezielles Steingut. Diese muss eingerichtet und angepasst werden. Man geht davon aus, dass die Anlage in höchstens neun Tagen einsatzbereit ist. Wir definieren als Ereignisse

A: Es dauert mehr als 6 Tage, bis die Anlage einsatzbereit ist. B: Es dauert weniger als 8 Tage, bis die Anlage einsatzbereit ist.

1. Beschreiben Sie das Komplement zu A.
2. Beschreiben Sie die Schnittmenge zwischen A und B.
3. Sind A und B disjunkt?
4. Zeigen Sie, dass $(A \cap B) \cup (\bar{A} \cap B) = B$ ist.

Aufgabe 1.2 - Einrichtung der Mischanlage

Wir beziehen uns auf Aufgabe 1.1 Die Werksleitung vermutet für die Zahl der Tage, die benötigt werden, um die Anlage einzurichten, die Wahrscheinlichkeiten, wie in der folgenden Tabelle angegeben:

Zahl der Tage	5	6	7	8	9
Wahrscheinlichkeit	0,05	0,25	0,35	0,25	0,10

1. Geben Sie die Wahrscheinlichkeiten für A und B an.
2. Geben Sie die Wahrscheinlichkeit für die Schnittmenge von A und B an.
3. Geben Sie die Wahrscheinlichkeit für die Vereinigungsmenge von A und B an.
4. Jeder unproduktive Tag kostet die Firma 2000 Euro. Mit welcher Wahrscheinlichkeit muss die Firma mit Kosten von höchstens 12.000 Euro rechnen?

Aufgabe 1.3 Zustelldienst

Ein Zustellungsdienst beschäftigt Festangestellte und freie Mitarbeiter. 64% der Mitarbeiter sind fest angestellt. Eine Qualitätsanalyse ergab, dass 10% aller Zustellungen beanstandet wurden. Die Wahrscheinlichkeit, dass eine Sendung von einem festangestellten Mitarbeiter ausgeliefert wurde und beanstandet wurde, beträgt 6%.

Berta erhält eine Sendung. Mit welcher Wahrscheinlichkeit

1. wird die Sendung beanstandet?
2. stammt die Sendung von einem freien Mitarbeiter?
3. wird die Sendung beanstandet oder stammt von einem Festangestellten?
4. wird die Sendung nicht beanstandet oder stammt nicht von einem Festangestellten?
5. wird die Sendung beanstandet oder stammt nicht von einem Festangestellten?
6. wird die Sendung beanstandet oder nicht beanstandet?
7. wird die Sendung beanstandet, stammt aber nicht von einem Festangestellten?

Aufgabe 1.4 - 2x Würfeln

Sie würfeln zweimal.

1. Geben Sie die Ergebnismenge dieses Zufallsvorgangs an. Zweckmäßig ist eine matrixähnliche Anordnung.
2. Mit welcher Wahrscheinlichkeit erhalten Sie
 1. beim ersten Wurf 1 und beim zweiten Wurf 5?
 2. einen Pasch (2x die gleiche Augenzahl)?
 3. 1 oder 5?
 4. die Augenzahl 8?
 5. mindestens die Augenzahl 7?

Aufgabe 1.5 - Münze 4x werfen

Eine Münze wird viermal geworfen. Es ist definiert: Z: Zahl liegt oben. K: Kopf liegt oben.

1. Stellen Sie die 16-elementige Ergebnismenge zusammen.

2. Geben Sie ein Beispiel für ein Ergebnis, ein Elementarereignis, ein zusammengesetztes Ereignis.

3. Es sind die Ereignisse definiert:

A: Es treten zuerst zweimal Kopf, dann zweimal Zahl auf

B: Es tritt höchstens zweimal Kopf auf

C: Es tritt mindestens drei mal Zahl auf

D: Es tritt einmal Kopf auf

Ermitteln Sie die Wahrscheinlichkeit, dass

1. A
2. D
3. nicht D
4. B und C
5. B oder C
6. A und C
7. nicht C und nicht D
8. nicht Kopf und nicht mindestens 3 mal Zahl
9. B ohne A
10. mindestens einmal Zahl

eintritt.

Aufgabe 1.6 - Aktiengewinne

Die Wertpapierabteilung einer Bank verwendet einen neuen speziellen Index zur Bewertung der zukünftigen Ertragsstärke eines Unternehmens. Eine erste Analyse ihrer Aktienportefeuilles hat ergeben, dass 75% der Aktien, deren Unternehmen als ertragsstark eingestuft worden waren, Kursgewinne einfahren konnten. Es wurden aber auch mit 30% der Aktien als ertragsschwach beurteilten Unternehmen Gewinne erzielt. Zur Vermeidung von Risiken setzten sich die Wertpapierfonds aus 80% Aktien als ertragsstark und 20% Aktien als ertragsschwach beurteilten Unternehmen zusammen.

1. Mit welcher Wahrscheinlichkeit kann von einer Aktie ein Kursgewinn erwartet werden?
2. Wieviel Prozent der Aktien mit Kursverlusten stammten tatsächlich von als ertragsschwach beurteilten Unternehmen?

Aufgabe 1.7 - Heulomat

Die Auto-Alarmanlage Heulomat heult erfahrungsgemäß bei 90% der Autoknacker, die sich am Auto zu schaffen machen. Leider heult sie auch bei 60% aller harmlosen Kollisionen, beispielsweise mit Spaziergängern. Man vermutet, dass insgesamt 80% aller Erschütterungen eines Autos harmlos sind.

1. In wie viel Prozent aller Fälle heult die Anlage **berechtigerweise**?
2. Wie groß ist die Wahrscheinlichkeit, dass die Anlage bei irgendeiner Erschütterung nicht heult?

Aufgabe 1.8 - Zugverspätung

Das Eisenbahnsystem des Staates Mobilia weist im Prinzip zwei Typen von Zügen auf: Den überregionalen Schnellzug „Hypercity“ und die langsamere Regionalbahn „Bummelzug“. Der Anteil der Hypercities am Fahrzeugbestand beträgt 20%. Man hat herausgefunden, dass 70% aller Hypercities verspätet sind, wogegen 80% aller Bummelzüge pünktlich ankommen.

Sie stehen am Bahnhof von Capitalis, der Hauptstadt von Mobilia, und sehen dem Treiben an den Bahnsteigen zu. Eine Lautsprecherdurchsage verkündet: „Der Zug nach Metropolis fährt verspätet ein“.

1. Mit welcher Wahrscheinlichkeit handelt es sich um
 1. einen Bummelzug?
 2. einen Hypercity ?
2. Wie groß ist der Anteil der Züge mit Verspätung?

Aufgabe 1.9 - Wand verkratzen mit Mülltonnen

In einem Mietshaus wird Dienstags die Mülltonne entleert. Bei 30% der Leerungen stellt Herr Löhlein die Mülltonne raus, bei 20% der Leerungen Frau Susemihl und bei 50% aller Leerungen Herr Feinbein. Eines Tages stellt der Vermieter fest, dass die Wand im Flur verschrammt ist. Er weiß, dass Herr Löhlein beim Mülltonne Tragen mit einer Wahrscheinlichkeit von 7%, Frau Susemihl mit einer Wahrscheinlichkeit von 8% und Herr Feinbein mit einer Wahrscheinlichkeit von 5% mit der Tonne an der Wand entlang kratzen.

1. Welcher Bewohner ist am „verdächtigsten“?
2. Mit welcher Wahrscheinlichkeit wird nächsten Dienstag die Wand verkratzt?
3. Nach jeder Schramme lässt der Vermieter die Wand weißen. Reicht etwa ein Anstrich pro Jahr?

Aufgabe 1.10 - Kaffeetassen

Frau Ahorn, Frau Behorn und Frau Zehorn bestellen nacheinander (in der Reihenfolge der Nennung) im Café Linde Kaffee. Zur Zeit sind noch 24 graue Tassen und 12 rosa Tassen heil. Die Tassen werden in der Reihenfolge der Bestellung zufällig ausgegeben.

1. Wie groß ist die Wahrscheinlichkeit, daß Frau Ahorn eine graue, Frau Behorn und Frau Zehorn eine rosa Tasse erhalten?
2. Wie groß ist die Wahrscheinlichkeit, daß Frau Zehorn eine rosa Tasse erhält?
3. Wie groß ist die Wahrscheinlichkeit, daß mindestens eine Kundin eine graue Tasse erhält?
4. Wie groß ist die Wahrscheinlichkeit, daß genau eine Kundin eine rosa Tasse erhält?
5. Es betreten 10 Kundinnen das Café. Wie groß ist die Wahrscheinlichkeit, daß mindestens 9 Kundinnen eine rosa Tasse erhalten?
(Ansatz genügt)

Aufgabe 1.11 Kondensatoren

Einem Fertigungslos von 500 Kondensatoren werden fünf Kondensatoren zu Prüfzwecken entnommen. Aufgrund einer ungenauen Wicklung sind 100 schadhafte Kondensatoren im Fertigungslos. Mit welcher Wahrscheinlichkeit taucht kein einziger dieser schadhafte Kondensatoren in der Probe auf?

Aufgabe 1.12 - Schraubensortiment

Einem Heimwerkemarkt werden Schachteln mit Schraubensortimenten geliefert, die jeweils 30 kleine Schrauben, 20 mittlere Schrauben und 10 große Schrauben enthalten. Zu Kontrollzwecken werden den Schachteln Schrauben entnommen.

1. Es wird 3 Schachteln jeweils eine Schraube entnommen. Wie groß ist die Wahrscheinlichkeit,
 1. dass erst eine kleine, dann eine große, dann eine mittlere Schraube resultiert?
 2. dass mindestens eine große Schraube resultiert?
2. Es werden einer Schachte drei Schrauben (o. Z.) entnommen. Wie groß ist die Wahrscheinlichkeit, dass nur kleine und mittlere Schrauben gezogen werden?

Kapitel 2

Aufgabe 2.1 Münze 3x werfen

Eine Münze wird dreimal geworfen.

1. Geben Sie die acht-elementige Ergebnismenge für den Zufallsvorgang: „Eine Münze wird dreimal geworfen“ an (K: Kopf; Z.: Zahl).
2. Definiert ist die Zufallsvariable X: Anzahl von Kopf bei drei Würfeln.
 1. Bestimmen Sie die Wahrscheinlichkeitsfunktion von X
 2. Berechnen Sie den Erwartungswert und die Varianz von X
3. Der Zufallsvorgang ist die Grundlage für ein Glücksspiel. Eine Person zahlt einen Einsatz von 1 Euro. Sie wirft dreimal eine Münze. Für jeden Kopf erhält sie 60 Cents. Es sei die Zufallsvariable Y der Nettogewinn.
 1. Geben Sie die Wahrscheinlichkeitsfunktion von Y an. Bestimmen Sie daraus $E(Y)$ und $VAR(Y)$.
 2. Geben Sie Y in Abhängigkeit von X an.
 3. Überlegen Sie, ob Y eine lineare Transformation von X ist.
 4. Berechnen Sie gegebenenfalls die Parameter von Y mit Hilfe dieser Erkenntnis.

Aufgabe 2.2 - Urne mit Kugeln

In einer Urne befinden sich 3 rote und 7 blaue Kugeln. Der Urne werden 4 Kugeln ohne Zurücklegen entnommen.

1. Mit welcher Wahrscheinlichkeit erhalten Sie
 1. keine rote Kugel?
 2. mindestens 1 rote Kugel?
 3. vier rote Kugeln?
2. Es sei definiert X: Zahl der roten Kugeln bei $n=4$.
 1. Geben Sie für X die Wahrscheinlichkeitstabelle und die Verteilungsfunktion an.
 2. Tragen Sie die Verteilungsfunktion in ein Diagramm ein. Hinweis: Es genügt, wenn Sie für die Ordinate im Nenner 210 stehen lassen.
 3. Geben Sie Erwartungswert und Varianz von X an.

Aufgabe 2.3 - Buchladen

Eine Buchhandlung steht vor der Wahl, ein hochwertiges und sehr teures Faksimile einer mittelalterlichen Handschrift anzubieten. Die Marketingexperten eines beauftragten Instituts vermuten für die Verkaufszahlen X folgende Wahrscheinlichkeiten:

Verkaufszahl x	0	1	2	3	4	5	mehr als 5
Wahrscheinlichkeit	0,2	0,3	0,2	0,1	0,1	0,1	0

1. Zeichnen Sie die Verteilungsfunktion.
2. Bestimmen Sie die Wahrscheinlichkeit, dass
 1. höchstens ein Buch
 2. weniger als zwei Bücher
 3. mindestens vier Bücher
 4. mehr als ein, aber höchstens vier Bücherverkauft werden.
1. Bestimmen Sie die durchschnittliche Zahl von Bücher, die eine Buchhandlung verkaufen könnte, und die Varianz.

Aufgabe 2.4 - Bäckerei

Die Bäckerei Körnchen hat festgestellt, dass sich die Zahl der täglich verkauften Mischbrote annähernd durch die Zufallsvariable X (in 100) mit einer Dichtefunktion

$$f(x) = \begin{cases} ax & \text{für } 0 \leq x \leq 6 \\ 0 & \text{sonst} \end{cases}$$

beschreiben lässt.

1. An wie viel Prozent der Tage können höchstens 400 Brote verkauft werden?
2. An wie viel Prozent der Tage können mindestens 500 Brote verkauft werden?
3. An wie viel Prozent der Tage können zwischen 400 und 500 Brote verkauft werden?
4. An wie viel Prozent der Tage können genau 600 Brote verkauft werden?
5. Bestimmen Sie a so, dass f tatsächlich eine Dichtefunktion ist.
6. Bestimmen Sie analytisch Verteilungsfunktion, Erwartungswert und Varianz von X.
7. Geben Sie den Median der Verteilung an.
8. Wie viel Brote wurden mindestens an den 20% „besten“ Tagen verkauft?

Aufgabe 2.5 - 2x Würfeln

Sie würfeln zweimal. Es ist die Zufallsvariable Y definiert als Summe der Augenzahlen der beiden Würfe.

1. Geben Sie Wahrscheinlichkeitstabelle und Verteilungsfunktion von Y an. Erstellen Sie jeweils eine Grafik.
2. Geben Sie die Wahrscheinlichkeit an,
 1. dass die Summe der Augenzahlen genau 4 beträgt.
 2. dass die Summe der Augenzahlen genau 2,5 beträgt.
 3. dass die Summe der Augenzahlen mindestens 4 beträgt.
 4. dass die Summe der Augenzahlen mehr als 4 beträgt.
 5. dass die Summe der Augenzahlen mehr als 9,5 beträgt.
 6. dass die Summe der Augenzahlen höchstens 3 beträgt.
 7. dass die Summe der Augenzahlen mindestens 4 und höchstens 10 beträgt.
 8. dass die Summe der Augenzahlen mindestens 4 oder höchstens 10 beträgt.
 9. dass Y mehr als 6 und weniger als 8 beträgt.
3. Bestimmen Sie Erwartungswert und Varianz von Y

Aufgabe 2.6 - Gemeinsame Wahrscheinlichkeiten

Die gemeinsamen Wahrscheinlichkeiten der diskreten Zufallsvariablen X und y sind in der folgenden Wahrscheinlichkeitstabelle zusammengefasst:

X \ Y	-2	-1	0	1	$f_{X^*}(x_i)$
0	0,05	0,05	0,05	0,1	
1	0	0,1	0,2	0,05	
2	0	0	0,2	0,1	
3	0	0	0	0,1	
$f_Y(y_j)$					

1. Bestimmen Sie Verteilung, Erwartungswert und Varianz von X und Y.
2. Überprüfen Sie, ob X und Y stochastisch unabhängig sind.
3. Ermitteln Sie den Korrelationskoeffizienten von X und Y.

Aufgabe 2.7 - Rendite zweier Aktien

Die Studentin Berta möchte das Geld, das sie durch Programmieraufträge verdient hat, in Aktien anlegen. Ihr erscheinen die Newcomer Scheffel und Raff am aussichtsreichsten. Sie hat die Wahrscheinlichkeiten für die Renditen (in Croetos), die die beiden Aktien gemeinsam abwerfen, in einer Renditetabelle zusammengefasst:

Scheffel	Raff	Wahrscheinlichkeit
X	Y	f_{XY}
0	0	0,1
0	10	0,1

50	10	0,2
50	30	0,1
100	30	0,2
100	40	0,3

1. Geben Sie die gemeinsame Wahrscheinlichkeitstabelle von X und Y an.
2. Ermitteln Sie die durchschnittliche Rendite einer Aktie und ihre Varianz.
3. Ermitteln Sie den Korrelationskoeffizienten zwischen den Renditen.
4. Berta zahlt ihrem Anlageverwalter jährlich einmal 10 € und dann von der Rendite 1%. Wieviel muss sie ihrem Anlageverwalter jährlich im Durchschnitt zahlen, wenn sie Scheffel und Raff kaufen würde?

Kapitel 3

Aufgabe 3.1 - Abnahmekontrolle von Elektronik

Bei einer sehr großen Lieferung von hochwertigen elektronischen Bauteilen wird ein Ausschussanteil von 5% als akzeptabel angesehen. Bei der Abnahmekontrolle werden 15 Stück zufällig entnommen. Falls höchstens ein fehlerhaftes Stück auftritt, wird die Lieferung angenommen.

1. Bestimmen Sie die Wahrscheinlichkeit,
 1. dass die Lieferung angenommen wird, wenn tatsächlich 5% Ausschuss vorliegen.
 2. dass die Lieferung irrtümlicherweise abgelehnt wird, wenn tatsächlich 3% Ausschuss vorliegen.
 3. dass die Lieferung irrtümlicherweise angenommen wird, wenn tatsächlich 10% Ausschuss vorliegen.
2. Wie groß muss die Stichprobe mindestens sein, damit die Wahrscheinlichkeit für eine irrtümliche Annahme der Lieferung bei 10% Ausschuss höchstens 10% beträgt? Verwenden Sie dazu die [Binomialverteilungstabelle](#).
3. Oft vermeidet man die Abnahmeregel, dass eine Lieferung nur angenommen wird, wenn kein fehlerhaftes Stück auftritt, weil man diese Regel für zu streng hält. Beurteilen Sie diese Ansicht.

Aufgabe 3.2 - Bank nach 18 Uhr

Die Zahl der Kunden, die nach 18 Uhr während einer Stunde einen Bankschalter in einer Bankfiliale aufsuchen, ist poissonverteilt mit $\lambda = 10$.

1. Wie viele Kunden suchen stündlich im Durchschnitt einen Bankschalter auf?
2. Wie groß ist der Anteil der Stunden, in denen höchstens drei Kunden an einen Schalter kommen?
3. Wie groß ist der Anteil der Stunden, in denen mindestens zwei Kunden an einen Schalter kommen?

Aufgabe 3.3 - LKW-Versicherung

Die Zahl der Versicherungsfälle, die einer gewerblichen Haftpflichtversicherung durch einen LKW entstehen, ist annähernd poissonverteilt mit dem Parameter $\lambda = 2,5$.

1. Bei wie viel Prozent der LKWs muss die Versicherung in einem Jahr keinen Schadensersatz leisten?
2. Wie viel Prozent der LKWs verursachen mindestens drei Versicherungsleistungen?
3. Eine Firma betreibt für just in time Lieferungen drei LKWs. Verursacht keiner der LKWs Versicherungsleistungen, bekommt die Firma 2000 Euro gutgeschrieben, falls doch, ändert sich finanziell nichts für die Firma. Ist das Angebot der Versicherung Ihrer Meinung nach attraktiv für die Firma?

Aufgabe 3.4 - Batteriefunktion

Für die Tauglichkeitsprüfung eines MP3-Players wurde geprüft, wie lange man ihn mit einem Batteriensatz spielen kann. Es stellte sich heraus, dass die Funktionsdauer eines Batteriensatzes annähernd normalverteilt ist mit dem Erwartungswert von 200 Minuten und einer Standardabweichung von 20 Minuten.

1. Bestimmen Sie die Wahrscheinlichkeit, dass ein MP3-Player mit einem Batteriensatz höchstens drei Stunden aushält.
2. Wie viel Prozent der MP3-Player schaffen mindestens 150 Minuten?
3. Mit welcher Wahrscheinlichkeit spielt ein MP3-Player zwischen zweieinhalb und dreieinhalb Stunden?
4. Bestimmen Sie d derart, dass der Anteil der MP3-Player, die zwischen $\mu - d$ und $\mu + d$ aushalten, 90% beträgt.

Aufgabe 3.5 - Küchenschaben

Eine Diplomarbeit über Küchenschaben hat ergeben, dass die Länge von Küchenschaben in einer bestimmten Altbauwohnung normalverteilt ist mit dem Erwartungswert 3 cm und der Varianz 4 cm^2 . In der Nacht wird eine Schabe zufällig eingefangen.

Bestimmen Sie die Wahrscheinlichkeit, dass diese Schabe

1. mindestens 5 cm
2. zwischen 2 und 5 cm
3. höchstens 1 cm
4. höchstens 2 oder mindestens 4 cm

lang ist.

Welche Mindestgröße haben die 10% größten Schaben?

Aufgabe 3.6 - Galapagos

Bei einer umfassenden Bestandsaufnahme von Großechsen auf einer Galapagosinsel stellte sich heraus, dass das Gewicht X dieser Echten annähernd normalverteilt ist. 15,87% der Echten wogen mehr als 120 kg. $x(0,33)$ betrug 75.

1. Tragen Sie die Angaben in die Grafik ein, wobei die Eintragungen nicht exakt maßstabsgetreu sein müssen.
2. Wieviel wogen die Echten im Durchschnitt?
3. Wieviel betrug die durchschnittliche quadratische Abweichung der Gewichte vom Mittel?



Aufgabe 3.7 - Nähfehler

Es ist bekannt, dass in einem Unternehmen, das Unterwäsche produziert, der Anteil von Spitzen-Damenunterhemden mit Nähfehlern etwa 10% beträgt. Der tägliche Output ist sehr groß. Es werden während eines Tages für die Warenkontrolle $n=200$ Hemdchen zufällig ausgewählt.

1. Bestimmen Sie die exakte Wahrscheinlichkeit, dass mindestens 15 Hemdchen Mängel aufweisen (nur Ansatz).
2. Berechnen Sie, falls möglich, die obige Wahrscheinlichkeit näherungsweise.
3. Mit welcher Wahrscheinlichkeit erhält man bei der Qualitätskontrolle mindestens 20 und höchstens 30 Hemdchen mit Fehlern?
4. Mit welcher Wahrscheinlichkeit erhält man bei der Qualitätskontrolle genau 20 fehlerhafte Hemdchen?
5. Ist es wahrscheinlicher, 19 bis 21 oder 23 bis 25 fehlerhafte Hemdchen zu erhalten?

Kapitel 4

Kapitel 5

Aufgabe 5.1 - Hotelsterne

Eine Reiseveranstalter hat 9 Kunden nach ihrer Zufriedenheit mit dem Hotel befragt, das sie im letzten Urlaub hatten.

Kunde Nr.	Sterne des Hotels	Note des Kunden
1	*	3
2	***	2
3	**	2
4	**	4
5	***	1
6	**	1
7	**	3
8	****	1
9	*	4

Ermitteln Sie den Rangkorrelationskoeffizienten der Sterne mit der Zufriedenheit

Aufgabe 5.2 - Solaranlagen

Eine Heizungsfirma hat in den letzten 8 Monaten jeweils x mal in der regionalen Tageszeitung inseriert. Sie konnte in diesen Monaten jeweils y viele Solaranlagen verkaufen.

Es ergab sich

Monat	i	1	2	3	4	5	6	7	8
Inserate	x	0	2	2	4	4	6	6	8
Solaranlagen	y	6	6	8	8	12	8	16	16

1. Tragen Sie die Wertepaare in einem Streudiagramm ab.
2. Ermitteln Sie die Regressionsgerade $y = a + bx$ und tragen Sie sie in das Diagramm ein.
3. Berechnen Sie die geschätzten Werte \hat{y} und die Residuen.
4. Berechnen Sie das Bestimmtheitsmaß.
5. Ermitteln Sie die Varianzen von y , \hat{y} und der Residuen. Zeigen Sie, dass die Streuungszerlegung hier gilt und ermitteln Sie das Bestimmtheitsmaß als Anteil der durch \hat{y} erklärten Streuung an der Gesamtstreuung von y .

Aufgabe 5.3 - Fair-Trade-Tee

Gegeben ist für die Jahre 1998 bis 2003 die Zahl der in der EU verkauften Tonnen Tee im fairen Handel.

Jahr	Zeitpunkt x	Menge y
1998	1	612
1999	2	842
2000	3	890
2001	4	1004
2002	5	1154
2003	6	1414

1. Ermitteln Sie eine Regressionsgerade, die die Entwicklung des Verkaufs im Lauf der Jahre beschreibt.
2. Berechnen Sie das Bestimmtheitsmaß.

Kapitel 6

Kapitel 7

Aufgabe 7.1 - Tarifsysteem

Eine Analyse der Kundenzufriedenheit eines großen Verkehrsbetriebes gab Anlass zu der Befürchtung, dass 75% der Fahrgäste das Tarifsysteem nicht verstanden hätten.

1. 75% der Kunden haben das Tarifsysteem nicht verstanden. Es wurden in einem zentral gelegenen U-Bahnhof zufällig 10 Personen befragt.
 1. Mit welcher Wahrscheinlichkeit hat jeder die Tarifordnung verstanden?
 2. Mit welcher Wahrscheinlichkeit haben genau 8 Personen die Tarifordnung verstanden?
 3. Mit welcher Wahrscheinlichkeit haben mindestens 2 Personen die Tarifordnung nicht verstanden?
 4. Mit welcher Wahrscheinlichkeit haben an zwei aufeinanderfolgenden Tagen jeweils mindestens zwei Personen das Tarifsysteem nicht verstanden, wenn die Befragungen stochastisch unabhängig waren.
2. Es wurden 100 Personen befragt.
 1. Mit welcher Wahrscheinlichkeit haben genau 75 Personen die Tarifordnung nicht verstanden?
 2. Mit welcher Wahrscheinlichkeit haben höchstens 75 Personen die Tarifordnung nicht verstanden?
 3. Es haben 70 Kunden angegeben, das System nicht verstanden zu haben. Überprüfen Sie die Hypothese ($\alpha = 0,05$), dass mindestens 75% die Tarifordnung nicht verstanden haben.

Aufgabe 7.2 - Kaviar

Ein Delikatessengroßhandel erhält eine umfangreiche Lieferung von 50-g-Schalen Kaviar. Es ist bekannt, dass die Füllmenge des Kaviars normalverteilt ist. Der Lieferant versichert, dass sich in jeder Dose im Mittel mindestens 50 g Kaviar befinden. Es werden zu Prüfzwecken 6 Schälchen zufällig ausgewählt und geöffnet. Man erhält die Urliste

47 49 50 52 50 46

1. Prüfen Sie die Behauptung des Lieferanten ($\alpha = 0,1$).
2. Würde sich die Position des Lieferanten verschlechtern, wenn man ein Signifikanzniveau von 0,05 verwenden würde?

9.

Statistik auf dem Computer

[<< hoch zum Anfang vom Inhaltsverzeichnis](#)

[<< hoch zum Kapitel 9 vom Inhaltsverzeichnis](#)

Einfache Statistikprogramme

Statistische Berechnungen mit der Tabellenkalkulation von Open Office

Auf vielen Rechnern findet sich mittlerweile das Programm Open Office.

Dort können Sie innerhalb der Tabellenkalkulation viele statistische Funktionen ausführen.

Öffnen Sie dazu ein neues Dokument in Open Office und wählen Sie bei der Art des Dokumentes *Tabellenkalkulation*.

Geben Sie Ihre Werte beispielsweise in der ersten Spalte A ein.

Über *Einfügen, Funktion* können Sie verschiedene Berechnungen auswählen, die Sie in einem freien Feld durchführen lassen. Ihre Werteliste markieren Sie für die Berechnung. Das Ergebnis erhalten Sie nach anklicken des grünen Häckchens neben der Eingabezeile.

Probieren Sie das ganze mit folgender Werteliste aus. Man kann sie mittels Zwischenspeicher direkt in die Tabellenkalkulation übernehmen.

(Anmarkieren, mit Strg + C in den Zwischenspeicher holen, im ersten Feld der Tabellenkalkulation mit Strg + V wieder abladen.)

```
114,3
135,7
104,8
118,5
125,7
121,4
122,4
96,8
118,9
120
112,2
127,9
122,8
128,9
120,3
```

Versuchen Sie den Median, den Mittelwert, die Standardabweichung, die Varianz und den Maximalwert zu ermitteln. Dazu gehen Sie in ein leeres Feld am Ende ihrer Liste. Dann wählen Sie im Menü den Eintrag *Einfügen* und dort wieder *Funktionsliste*. Aus der Funktionsliste wählen Sie die statistischen Funktionen aus. Aus dem großen Angebot wählen Sie den *Median*.

Sie können auch direkt in die Eingabezeile folgendes eingeben:

```
=MEDIAN (A1:A15)
```

oder

```
=MITTELWERT (A1:A15)
```

Es müßten folgende Werte herauskommen:

Summe: 1790,6

Mittelwert 119,37

Median 120,3

Maximal 135,7

Standardabweichung 9,62

Varianz 92,6

Statistik mit Gnumeric

Gnumeric ist die Tabellenkalkulation unter Gnome. Es bietet bessere statistische Berechnungsmöglichkeiten als Excel. Siehe <http://de.wikipedia.org/wiki/Gnumeric> ☞ Siehe <http://www.gnome.org/projects/gnumeric/> ☞

Statistische Berechnungen mit der Programmiersprache Gambas

Auf vielen Linuxrechnern findet sich mittlerweile die einfach zu lernende Programmsprache Gambas.

Dort kann man viele statistische Funktionen nachvollziehen. Es gilt das alte Motto: Habe ich es noch nicht programmiert, dann habe ich es noch nicht verstanden.

Im Gambas Wikibook sollen nach und nach eine Reihe von Statistikfunktionen im Quelltext erklärt und verfügbar gemacht werden.

Siehe http://de.wikibooks.org/wiki/Gambas:_Statistik ☞

Komplexere, professionelle Statistik-Software

R (Windows, OS X, Linux)

R ist eine umfangreiche Statistiksoftware, genauer: eine *Programmierungsumgebung* für statistische Auswertungen. Im Funktionsumfang mit kommerziellen Softwarepaketen wie SPSS oder STATA durchaus vergleichbar (und stellenweise überlegen) werden eine Vielzahl statistischer Methoden und Routinen bereitgestellt. Der Programmaufbau mag zwar insbesondere für Anfänger etwas unübersichtlich sein, besticht jedoch durch zahlreiche Features: Vollständige Kontrolle über die Daten, Implementation einer grossen Anzahl an Analyse-Verfahren, flexible Graphikfähigkeiten, Systemunabhängigkeit, automatisierte Auswertungen, Schnittstellen zu vielen anderen Anwendungen und nicht zuletzt kostenlose Verfügbarkeit könnten dazu führen, dass R sich im professionellen Bereich zum neuen Standard entwickelt.



Tipp: Eine Einführung in R bietet das Wikibook [GNU R](#).

Einsteiger können zudem auf graphische Bedienoberflächen zurückgreifen.

Siehe:

- R
 - <http://www.r-project.org/>
 - http://de.wikibooks.org/wiki/GNU_R
- Graphische Bedienung (GUIs):
 - Jaguar: <http://stats.math.uni-augsburg.de/JGR/>
 - R Commander: <http://socserv.mcmaster.ca/jfox/Misc/Rcmdr/>
 - Das Statistiklabor: <http://www.statistiklabor.de>

SPSS (Windows, OS X, Linux (nur Server-Variante))

Insbesondere in den Sozial- und Verhaltenswissenschaften findet die kommerzielle Software SPSS grossen Zuspruch. Hauptvorteile sind einfache Bedienbarkeit für Anwender, die die Steuerung mit Menüs und Maus favorisieren. Tabellen und Graphiken sind für die Weiterverwendung in Office-Anwendungen optimiert und können nachträglich formatiert werden. Zudem gibt es Ergänzungspakete, die den Prozess der Datenerfassung erleichtern. Integriert sind die meisten gebräuchlichen Standardverfahren uni- und multivariater Statistik. Spezielle Anwendungen werden als Zusatzpakete vertrieben.

- [Wikipedia über SPSS](#)
- <http://www.spss.com/de/>

STATA (Windows, OS X, Linux)

STATA ist eine Statistiksoftware, die bislang hauptsächlich im anglo-amerikanischen Raum Verbreitung gefunden hat. Wenngleich STATA ebenfalls über eine gut strukturierte Menübedienung verfügt, besticht die Software vor allem durch ihre an [BASIC](#) erinnernde, relativ einfach zu erlernende Befehlssyntax und eine umfängliche, übersichtliche Integration weitreichender statistischer Verfahren.

- [Wikipedia über STATA](#)
- <http://www.stata.com>

10.

Literaturliste

[<< hoch zum Anfang vom Inhaltsverzeichnis](#)

[<< hoch zum Kapitel 10 vom Inhaltsverzeichnis](#)

- Grundlagen der Statistik, Bd.1, Beschreibende Verfahren von Jochen Schwarze, Verlag Neue Wirtschafts-Briefe (Januar 2001)
- Grundlagen der Statistik, Bd.2, Wahrscheinlichkeitsrechnung und induktive Statistik von Jochen Schwarze, Verlag Neue Wirtschafts-Briefe (Oktober 2001)
- Aufgabensammlung zur Statistik von Jochen Schwarze Verlag Neue Wirtschafts-Briefe (Januar 2002)
- Wahrscheinlichkeitsrechnung und schließende Statistik . Praxisorientierte Einführung. Mit Aufgaben und Lösungen von Günther Bourier
- Statistik-Übungen von Günther Bourier
- Beschreibende Statistik von Günther Bourier
- Stochastik Leistungskurs Lambacher Schweizer, Klett Verlag
 - Didaktisch gut gemacht, mit sehr vielen Übungsaufgaben

- Stochastik Leistungskurs Lambacher Schweizer, Klett Verlag Lösungsheft
- Medizinische Statistik, Von Herbert Immich , Schattauer Verlag

Leider nur noch antiquarisch, didaktisch hervorragend gemacht, mit sehr vielen praktischen medizinischen Beispielen.

[FISZ89] Marek Fisz, *Wahrscheinlichkeitsrechnung und mathematische Statistik*, VEB Deutscher Verlag der Wissenschaften, Berlin 1989, [ISBN 3-326-00079-0](#)

Für Fortgeschrittene. Wie ich finde, ein umfassendes Lehrbuch mit einer guten Einführung in die Wahrscheinlichkeitsrechnung und vielen guten Beispielen. Die Testverfahren sind meiner Meinung nach gut, aber sehr theoretisch, beschrieben.

Helge Nordmann, *Statistik und Wahrscheinlichkeit – leicht gemacht*, BooksOnDemand, Norderstedt 2007, [ISBN 978-3-8334-9977-7](#)

Ein Lehr- und Nachschlagewerk. Zu jedem Thema/Kapitel finden sich Übungsaufgaben mit vollständigem Lösungsweg.

11.	Tabellen	<< hoch zum Anfang vom Inhaltsverzeichnis << hoch zum Kapitel 11 vom Inhaltsverzeichnis
------------	-----------------	--

11.1.	Verteilungstabellen	<< hoch zum Anfang vom Inhaltsverzeichnis << hoch zum Kapitel 11 vom Inhaltsverzeichnis
--------------	----------------------------	--

11.1.1.	Normalverteilung	<< hoch zum Anfang vom Inhaltsverzeichnis << hoch zum Kapitel 11 vom Inhaltsverzeichnis
----------------	-------------------------	--

Erläuterungen zur Normalverteilungstabelle:

Die fettgedruckten Werte sind die Ausprägungen z der standardnormalverteilten Zufallsvariablen Z , die vierstelligen Ziffern in der Tabelle selbst stellen die zu z gehörigen Verteilungswerte dar, wobei "0," weggelassen wurde.

Es ist beispielsweise $P(Z \leq 2,51) = 0,9940$.

z	0,00	0,01	0,02	0,03	0,04	0,05	0,06	0,07	0,08	0,09
0,0	5000	5040	5080	5120	5160	5199	5239	5279	5319	5359
0,1	5398	5438	5478	5517	5557	5596	5636	5675	5714	5753
0,2	5793	5832	5871	5910	5948	5987	6026	6064	6103	6141
0,3	6179	6217	6255	6293	6331	6368	6406	6443	6480	6517
0,4	6554	6591	6628	6664	6700	6736	6772	6808	6844	6879
0,5	6915	6950	6985	7019	7054	7088	7123	7157	7190	7224
0,6	7257	7291	7324	7357	7389	7422	7454	7486	7517	7549

0,7	7580	7611	7642	7673	7704	7734	7764	7794	7823	7852
0,8	7881	7910	7939	7967	7995	8023	8051	8078	8106	8133
0,9	8159	8186	8212	8238	8264	8289	8315	8340	8365	8389
1,0	8413	8438	8461	8485	8508	8531	8554	8577	8599	8621
1,1	8643	8665	8686	8708	8729	8749	8770	8790	8810	8830
1,2	8849	8869	8888	8907	8925	8944	8962	8980	8997	9015
1,3	9032	9049	9066	9082	9099	9115	9131	9147	9162	9177
1,4	9192	9207	9222	9236	9251	9265	9279	9292	9306	9319
1,5	9332	9345	9357	9370	9382	9394	9406	9418	9429	9441
1,6	9452	9463	9474	9484	9495	9505	9515	9525	9535	9545
1,7	9554	9564	9573	9582	9591	9599	9608	9616	9625	9633
1,8	9641	9649	9656	9664	9671	9678	9686	9693	9699	9706
1,9	9713	9719	9726	9732	9738	9744	9750	9756	9761	9767
2,0	9772	9778	9783	9788	9793	9798	9803	9808	9812	9817
2,1	9821	9826	9830	9834	9838	9842	9846	9850	9854	9857
2,2	9861	9864	9868	9871	9875	9878	9881	9884	9887	9890
2,3	9893	9896	9898	9901	9904	9906	9909	9911	9913	9916
2,4	9918	9920	9922	9925	9927	9929	9931	9932	9934	9936
2,5	9938	9940	9941	9943	9945	9946	9948	9949	9951	9952
2,6	9953	9955	9956	9957	9959	9960	9961	9962	9963	9964
2,7	9965	9966	9967	9968	9969	9970	9971	9972	9973	9974
2,8	9974	9975	9976	9977	9977	9978	9979	9979	9980	9981

2,9	9981	9982	9982	9983	9984	9984	9985	9985	9986	9986
3,0	9987	9987	9987	9988	9988	9989	9989	9989	9990	9990
3,1	9990	9991	9991	9991	9992	9992	9992	9992	9993	9993
3,2	9993	9993	9994	9994	9994	9994	9994	9995	9995	9995
3,3	9995	9995	9995	9996	9996	9996	9996	9996	9996	9997
3,4	9997	9997	9997	9997	9997	9997	9997	9997	9997	9998

11.1.2.

Chi-Quadrat-Verteilung

[<< hoch zum Anfang vom Inhaltsverzeichnis](#)

[<< hoch zum Kapitel 11 vom Inhaltsverzeichnis](#)

Quantile der Chi-Quadrat-Verteilung nach ausgewählten Wahrscheinlichkeiten p und Freiheitsgraden

	Wahrscheinlichkeit p										
Freiheitsgrade	0,005	0,01	0,025	0,05	0,1	0,5	0,9	0,95	0,975	0,99	0,995
1	0,00	0,00	0,00	0,00	0,02	0,45	2,71	3,84	5,02	6,63	7,88
2	0,01	0,02	0,05	0,10	0,21	1,39	4,61	5,99	7,38	9,21	10,60
3	0,07	0,11	0,22	0,35	0,58	2,37	6,25	7,81	9,35	11,34	12,84
4	0,21	0,30	0,48	0,71	1,06	3,36	7,78	9,49	11,14	13,28	14,86
5	0,41	0,55	0,83	1,15	1,61	4,35	9,24	11,07	12,83	15,09	16,75
6	0,68	0,87	1,24	1,64	2,20	5,35	10,64	12,59	14,45	16,81	18,55
7	0,99	1,24	1,69	2,17	2,83	6,35	12,02	14,07	16,01	18,48	20,28
8	1,34	1,65	2,18	2,73	3,49	7,34	13,36	15,51	17,53	20,09	21,95
9	1,73	2,09	2,70	3,33	4,17	8,34	14,68	16,92	19,02	21,67	23,59
10	2,16	2,56	3,25	3,94	4,87	9,34	15,99	18,31	20,48	23,21	25,19
p →	0,005	0,01	0,025	0,05	0,1	0,5	0,9	0,95	0,975	0,99	0,995
11	2,60	3,05	3,82	4,57	5,58	10,34	17,28	19,68	21,92	24,73	26,76

12	3,07	3,57	4,40	5,23	6,30	11,34	18,55	21,03	23,34	26,22	28,30
13	3,57	4,11	5,01	5,89	7,04	12,34	19,81	22,36	24,74	27,69	29,82
14	4,07	4,66	5,63	6,57	7,79	13,34	21,06	23,68	26,12	29,14	31,32
15	4,60	5,23	6,26	7,26	8,55	14,34	22,31	25,00	27,49	30,58	32,80
16	5,14	5,81	6,91	7,96	9,31	15,34	23,54	26,30	28,85	32,00	34,27
17	5,70	6,41	7,56	8,67	10,09	16,34	24,77	27,59	30,19	33,41	35,72
18	6,26	7,01	8,23	9,39	10,86	17,34	25,99	28,87	31,53	34,81	37,16
19	6,84	7,63	8,91	10,12	11,65	18,34	27,20	30,14	32,85	36,19	38,58
20	7,43	8,26	9,59	10,85	12,44	19,34	28,41	31,41	34,17	37,57	40,00
p →	0,005	0,01	0,025	0,05	0,1	0,5	0,9	0,95	0,975	0,99	0,995
21	8,03	8,90	10,28	11,59	13,24	20,34	29,62	32,67	35,48	38,93	41,40
22	8,64	9,54	10,98	12,34	14,04	21,34	30,81	33,92	36,78	40,29	42,80
23	9,26	10,20	11,69	13,09	14,85	22,34	32,01	35,17	38,08	41,64	44,18
24	9,89	10,86	12,40	13,85	15,66	23,34	33,20	36,42	39,36	42,98	45,56
25	10,52	11,52	13,12	14,61	16,47	24,34	34,38	37,65	40,65	44,31	46,93
26	11,16	12,20	13,84	15,38	17,29	25,34	35,56	38,89	41,92	45,64	48,29
27	11,81	12,88	14,57	16,15	18,11	26,34	36,74	40,11	43,19	46,96	49,65
28	12,46	13,56	15,31	16,93	18,94	27,34	37,92	41,34	44,46	48,28	50,99
29	13,12	14,26	16,05	17,71	19,77	28,34	39,09	42,56	45,72	49,59	52,34
30	13,79	14,95	16,79	18,49	20,60	29,34	40,26	43,77	46,98	50,89	53,67
p →	0,005	0,01	0,025	0,05	0,1	0,5	0,9	0,95	0,975	0,99	0,995
31	14,46	15,66	17,54	19,28	21,43	30,34	41,42	44,99	48,23	52,19	55,00
32	15,13	16,36	18,29	20,07	22,27	31,34	42,59	46,19	49,48	53,49	56,33
33	15,82	17,07	19,05	20,87	23,11	32,34	43,75	47,40	50,73	54,78	57,65
34	16,50	17,79	19,81	21,66	23,95	33,34	44,90	48,60	51,97	56,06	58,96

35	17,19	18,51	20,57	22,47	24,80	34,34	46,06	49,80	53,20	57,34	60,28
36	17,89	19,23	21,34	23,27	25,64	35,34	47,21	51,00	54,44	58,62	61,58
37	18,59	19,96	22,11	24,08	26,49	36,34	48,36	52,19	55,67	59,89	62,88
38	19,29	20,69	22,88	24,88	27,34	37,34	49,51	53,38	56,90	61,16	64,18
39	20,00	21,43	23,65	25,70	28,20	38,34	50,66	54,57	58,12	62,43	65,48
40	20,71	22,16	24,43	26,51	29,05	39,34	51,80	55,76	59,34	63,69	66,77

11.1.3.

F-Verteilung

[<< hoch zum Anfang vom Inhaltsverzeichnis](#)

[<< hoch zum Kapitel 11 vom Inhaltsverzeichnis](#)

Quantile der F-Verteilung für die Wahrscheinlichkeit 0,975 und k_1 und k_2 (1 bis 15)

Freiheitsgrade

$k_2 \rightarrow$	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
k_1															
1	647,79	38,51	17,44	12,22	10,01	8,81	8,07	7,57	7,21	6,94	6,72	6,55	6,41	6,30	6,20
2	799,48	39,00	16,04	10,65	8,43	7,26	6,54	6,06	5,71	5,46	5,26	5,10	4,97	4,86	4,77
3	864,15	39,17	15,44	9,98	7,76	6,60	5,89	5,42	5,08	4,83	4,63	4,47	4,35	4,24	4,15
4	899,60	39,25	15,10	9,60	7,39	6,23	5,52	5,05	4,72	4,47	4,28	4,12	4,00	3,89	3,80
5	921,83	39,30	14,88	9,36	7,15	5,99	5,29	4,82	4,48	4,24	4,04	3,89	3,77	3,66	3,58
6	937,11	39,33	14,73	9,20	6,98	5,82	5,12	4,65	4,32	4,07	3,88	3,73	3,60	3,50	3,41
7	948,20	39,36	14,62	9,07	6,85	5,70	4,99	4,53	4,20	3,95	3,76	3,61	3,48	3,38	3,29
8	956,64	39,37	14,54	8,98	6,76	5,60	4,90	4,43	4,10	3,85	3,66	3,51	3,39	3,29	3,20
9	963,28	39,39	14,47	8,90	6,68	5,52	4,82	4,36	4,03	3,78	3,59	3,44	3,31	3,21	3,12
10	968,63	39,40	14,42	8,84	6,62	5,46	4,76	4,30	3,96	3,72	3,53	3,37	3,25	3,15	3,06
$k_2 \rightarrow$	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
11	973,03	39,41	14,37	8,79	6,57	5,41	4,71	4,24	3,91	3,66	3,47	3,32	3,20	3,09	3,01

12	976,72	39,41	14,34	8,75	6,52	5,37	4,67	4,20	3,87	3,62	3,43	3,28	3,15	3,05	2,96
13	979,84	39,42	14,30	8,72	6,49	5,33	4,63	4,16	3,83	3,58	3,39	3,24	3,12	3,01	2,92
14	982,55	39,43	14,28	8,68	6,46	5,30	4,60	4,13	3,80	3,55	3,36	3,21	3,08	2,98	2,89
15	984,87	39,43	14,25	8,66	6,43	5,27	4,57	4,10	3,77	3,52	3,33	3,18	3,05	2,95	2,86
16	986,91	39,44	14,23	8,63	6,40	5,24	4,54	4,08	3,74	3,50	3,30	3,15	3,03	2,92	2,84
17	988,72	39,44	14,21	8,61	6,38	5,22	4,52	4,05	3,72	3,47	3,28	3,13	3,00	2,90	2,81
18	990,35	39,44	14,20	8,59	6,36	5,20	4,50	4,03	3,70	3,45	3,26	3,11	2,98	2,88	2,79
19	991,80	39,45	14,18	8,58	6,34	5,18	4,48	4,02	3,68	3,44	3,24	3,09	2,96	2,86	2,77
20	993,08	39,45	14,17	8,56	6,33	5,17	4,47	4,00	3,67	3,42	3,23	3,07	2,95	2,84	2,76
$k_2 \rightarrow$	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
21	994,30	39,45	14,16	8,55	6,31	5,15	4,45	3,98	3,65	3,40	3,21	3,06	2,93	2,83	2,74
22	995,35	39,45	14,14	8,53	6,30	5,14	4,44	3,97	3,64	3,39	3,20	3,04	2,92	2,81	2,73
23	996,34	39,45	14,13	8,52	6,29	5,13	4,43	3,96	3,63	3,38	3,18	3,03	2,91	2,80	2,71
24	997,27	39,46	14,12	8,51	6,28	5,12	4,41	3,95	3,61	3,37	3,17	3,02	2,89	2,79	2,70
25	998,09	39,46	14,12	8,50	6,27	5,11	4,40	3,94	3,60	3,35	3,16	3,01	2,88	2,78	2,69
26	998,84	39,46	14,11	8,49	6,26	5,10	4,39	3,93	3,59	3,34	3,15	3,00	2,87	2,77	2,68
27	999,54	39,46	14,10	8,48	6,25	5,09	4,39	3,92	3,58	3,34	3,14	2,99	2,86	2,76	2,67
28	1000,24	39,46	14,09	8,48	6,24	5,08	4,38	3,91	3,58	3,33	3,13	2,98	2,85	2,75	2,66
29	1000,82	39,46	14,09	8,47	6,23	5,07	4,37	3,90	3,57	3,32	3,13	2,97	2,85	2,74	2,65
30	1001,40	39,46	14,08	8,46	6,23	5,07	4,36	3,89	3,56	3,31	3,12	2,96	2,84	2,73	2,64

Quantile der F-Verteilung für die Wahrscheinlichkeit 0,975 und k_1 und k_2 (16 bis 30)

Freiheitsgrade															
$k_2 \rightarrow$	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30
k_1															
1	6,12	6,04	5,98	5,92	5,87	5,83	5,79	5,75	5,72	5,69	5,66	5,63	5,61	5,59	5,57

2	4,69	4,62	4,56	4,51	4,46	4,42	4,38	4,35	4,32	4,29	4,27	4,24	4,22	4,20	4,18
3	4,08	4,01	3,95	3,90	3,86	3,82	3,78	3,75	3,72	3,69	3,67	3,65	3,63	3,61	3,59
4	3,73	3,66	3,61	3,56	3,51	3,48	3,44	3,41	3,38	3,35	3,33	3,31	3,29	3,27	3,25
5	3,50	3,44	3,38	3,33	3,29	3,25	3,22	3,18	3,15	3,13	3,10	3,08	3,06	3,04	3,03
6	3,34	3,28	3,22	3,17	3,13	3,09	3,05	3,02	2,99	2,97	2,94	2,92	2,90	2,88	2,87
7	3,22	3,16	3,10	3,05	3,01	2,97	2,93	2,90	2,87	2,85	2,82	2,80	2,78	2,76	2,75
8	3,12	3,06	3,01	2,96	2,91	2,87	2,84	2,81	2,78	2,75	2,73	2,71	2,69	2,67	2,65
9	3,05	2,98	2,93	2,88	2,84	2,80	2,76	2,73	2,70	2,68	2,65	2,63	2,61	2,59	2,57
10	2,99	2,92	2,87	2,82	2,77	2,73	2,70	2,67	2,64	2,61	2,59	2,57	2,55	2,53	2,51
$k_2 \rightarrow$	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30
11	2,93	2,87	2,81	2,76	2,72	2,68	2,65	2,62	2,59	2,56	2,54	2,51	2,49	2,48	2,46
12	2,89	2,82	2,77	2,72	2,68	2,64	2,60	2,57	2,54	2,51	2,49	2,47	2,45	2,43	2,41
13	2,85	2,79	2,73	2,68	2,64	2,60	2,56	2,53	2,50	2,48	2,45	2,43	2,41	2,39	2,37
14	2,82	2,75	2,70	2,65	2,60	2,56	2,53	2,50	2,47	2,44	2,42	2,39	2,37	2,36	2,34
15	2,79	2,72	2,67	2,62	2,57	2,53	2,50	2,47	2,44	2,41	2,39	2,36	2,34	2,32	2,31
16	2,76	2,70	2,64	2,59	2,55	2,51	2,47	2,44	2,41	2,38	2,36	2,34	2,32	2,30	2,28
17	2,74	2,67	2,62	2,57	2,52	2,48	2,45	2,42	2,39	2,36	2,34	2,31	2,29	2,27	2,26
18	2,72	2,65	2,60	2,55	2,50	2,46	2,43	2,39	2,36	2,34	2,31	2,29	2,27	2,25	2,23
19	2,70	2,63	2,58	2,53	2,48	2,44	2,41	2,37	2,35	2,32	2,29	2,27	2,25	2,23	2,21
20	2,68	2,62	2,56	2,51	2,46	2,42	2,39	2,36	2,33	2,30	2,28	2,25	2,23	2,21	2,20
$k_2 \rightarrow$	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30
21	2,67	2,60	2,54	2,49	2,45	2,41	2,37	2,34	2,31	2,28	2,26	2,24	2,22	2,20	2,18
22	2,65	2,59	2,53	2,48	2,43	2,39	2,36	2,33	2,30	2,27	2,24	2,22	2,20	2,18	2,16
23	2,64	2,57	2,52	2,46	2,42	2,38	2,34	2,31	2,28	2,26	2,23	2,21	2,19	2,17	2,15
24	2,63	2,56	2,50	2,45	2,41	2,37	2,33	2,30	2,27	2,24	2,22	2,19	2,17	2,15	2,14

25	2,61	2,55	2,49	2,44	2,40	2,36	2,32	2,29	2,26	2,23	2,21	2,18	2,16	2,14	2,12
26	2,60	2,54	2,48	2,43	2,39	2,34	2,31	2,28	2,25	2,22	2,19	2,17	2,15	2,13	2,11
27	2,59	2,53	2,47	2,42	2,38	2,33	2,30	2,27	2,24	2,21	2,18	2,16	2,14	2,12	2,10
28	2,58	2,52	2,46	2,41	2,37	2,33	2,29	2,26	2,23	2,20	2,17	2,15	2,13	2,11	2,09
29	2,58	2,51	2,45	2,40	2,36	2,32	2,28	2,25	2,22	2,19	2,17	2,14	2,12	2,10	2,08
30	2,57	2,50	2,44	2,39	2,35	2,31	2,27	2,24	2,21	2,18	2,16	2,13	2,11	2,09	2,07

11.1.4.

t-Verteilung

[<< hoch zum Anfang vom Inhaltsverzeichnis](#)

[<< hoch zum Kapitel 11 vom Inhaltsverzeichnis](#)

Quantile der t-Verteilung nach ausgewählten
Wahrscheinlichkeiten p und Freiheitsgraden

Freiheitsgrade	Wahrscheinlichkeit p				
	0,900	0,950	0,975	0,990	0,995
1	3,078	6,314	12,706	31,821	63,656
2	1,886	2,920	4,303	6,965	9,925
3	1,638	2,353	3,182	4,541	5,841
4	1,533	2,132	2,776	3,747	4,604
5	1,476	2,015	2,571	3,365	4,032
6	1,440	1,943	2,447	3,143	3,707
7	1,415	1,895	2,365	2,998	3,499
8	1,397	1,860	2,306	2,896	3,355
9	1,383	1,833	2,262	2,821	3,250
10	1,372	1,812	2,228	2,764	3,169
p →	0,900	0,950	0,975	0,990	0,995
11	1,363	1,796	2,201	2,718	3,106

12	1,356	1,782	2,179	2,681	3,055
13	1,350	1,771	2,160	2,650	3,012
14	1,345	1,761	2,145	2,624	2,977
15	1,341	1,753	2,131	2,602	2,947
16	1,337	1,746	2,120	2,583	2,921
17	1,333	1,740	2,110	2,567	2,898
18	1,330	1,734	2,101	2,552	2,878
19	1,328	1,729	2,093	2,539	2,861
20	1,325	1,725	2,086	2,528	2,845
p →	0,900	0,950	0,975	0,990	0,995
21	1,323	1,721	2,080	2,518	2,831
22	1,321	1,717	2,074	2,508	2,819
23	1,319	1,714	2,069	2,500	2,807
24	1,318	1,711	2,064	2,492	2,797
25	1,316	1,708	2,060	2,485	2,787
26	1,315	1,706	2,056	2,479	2,779
27	1,314	1,703	2,052	2,473	2,771
28	1,313	1,701	2,048	2,467	2,763
29	1,311	1,699	2,045	2,462	2,756
30	1,310	1,697	2,042	2,457	2,750
1000	1,282	1,646	1,962	2,330	2,581

11.2.

Weltdaten des CIA World-Fact-Book

[<< hoch zum Anfang vom Inhaltsverzeichnis](#)

[<< hoch zum Kapitel 11 vom Inhaltsverzeichnis](#)

Weltdaten

Diese Daten sind dem [Worldfact-Book der CIA](#) [🔗](#) entnommen und daher public domain. Es handelt sich um Indikatoren für die Entwicklung von Ländern. Kleinere Inselarchipele und weitere kleinste Staaten wurden entfernt, wobei die Elimination keinem Kriterium gehorchte, sondern rein nach Gutdünken erfolgte.

Erläuterung der Variablen

Die Variablen sind folgendermaßen definiert:

Country	Popgrow	Birth	Death	Fert	DPPop	Pov	Inflation
Land	Wachstumsrate der Bevölkerung	Geburtenrate (Zahl der Geburten auf 1000 Einwohner)	Sterberate (Zahl der Todesfälle auf 1000 Einwohner)	Fruchtbarkeitsrate (Zahl der Geburten pro gebärfähiger Frau)	Bruttoinlandsprodukt pro Kopf in USD	Prozentsatz der Bevölkerung unter der Armutsgrenze	Inflationsrate
Unemploy		ExpMilDP	Explife	MedAge	AgeStrJ	AgeStrM	AgeStrO
Arbeitslosenquote		Anteil der Rüstungsausgaben am Bruttoinlandsprodukt	Lebenserwartung eines Neugeborenen	Median des Lebensalters	Anteil der bis 15jährigen an der Bevölkerung	Anteil der 15- bis 65jährigen an der Bevölkerung	Anteil der ab 65jährigen an der Bevölkerung
Internet	Phone	Populat	InfMortT		InfMortM		InfMortF
Zahl der Internetuser	Zahl der Festnetztelefone	Gesamtbevölkerung	Kindersterblichkeit (Zahl der Todesfälle von Kindern unter 1 Jahr pro 1000 Lebendgeborenen)		Kindersterblichkeit (Zahl der Todesfälle von Knaben unter 1 Jahr pro 1000 lebendgeborenen Knaben)		Kindersterblichkeit (Zahl der Todesfälle von Mädchen unter 1 Jahr pro 1000 lebendgeborenen Mädchen)

Weiterführenden Erläuterungen zu den Variablen finden sich im <https://www.cia.gov/library/publications/the-world-factbook/docs/notesanddefs.html> 📖 World Factbook

Daten

Country	Popgrow	Birth	Death	Fert	DPPop	Pov	Inflation	Unemploy	ExpMilDP	Explife	MedAge	AgeStrJ	AgeStrM	AgeStrO	Internet	Phone	Populat	InfMortT	InfMortM	InfMortF	
Afghanistan	0,0492	47,27	21,12	6,78	700	0,2300	0,0520		0,01	42,46	17,5	0,447	0,529	0,024	1000	33100	28.513.677	165,96	170,85	160,82	1000
Albania	0,0051	15,08	5,02	2,05	4500	0,3000	0,0240	0,1580	0,0149	77,06	28,2	0,264	0,653	0,083	30000	255000	3.544.808	22,31	23,01	21,54	30000
Algeria	0,0128	17,76	4,61	2,04	6000	0,2300	0,0350	0,2620	0,035	72,74	23,8	0,299	0,655	0,046	500000	2199600	32.129.324	32,16	36,06	28,06	500000
Andorra	0,0100	9,32	5,9	1,28	19000		0,0430	0,0000		83,5	39,7	0,15	0,716	0,134	24500	35000	69.865	4,05	4,39	3,7	24500
Angola	0,0193	45,14	25,86	6,33	1900	0,7000	0,7660		0,019	36,79	18,1	0,435	0,537	0,028	41000	96300	10.978.552	192,5	204,97	179,41	41000
Antigua and Barbuda	0,0060	17,7	5,55	2,27	11000		0,0040	0,1100		71,6	29,4	0,281	0,676	0,043	10000	38000	68.320	20,18	24,29	15,87	10000
Argentina	0,0102	17,19	7,57	2,24	11200	0,5170	0,1340	0,1730	0,013	75,7	29,2	0,259	0,636	0,105	4100000	8009400	39.144.753	15,66	17,6	13,63	4100000
Armenia	-0,0032	11,43	8,12	1,31	3500	0,5000	0,0480	0,2000	0,065	71,23	29,7	0,227	0,667	0,106	150000	562600	2.991.360	24,16	29,32	18,06	150000
Australia	0,0090	12,4	7,38	1,76	29000		0,0280	0,0600	0,028	80,26	36,3	0,201	0,672	0,128	9472000	10815000	19.913.144	4,76	5,16	4,34	9472000
Austria	0,0014	8,9	9,56	1,35	30000	0,0390	0,0140	0,0440	0,008	78,87	40	0,159	0,681	0,16	3730000	3881000	8.174.762	4,68	5,76	3,55	3730000
Azerbaijan	0,0052	19,81	9,76	2,39	3400	0,4900	0,0210	0,0110	0,026	63,25	27,3	0,27	0,652	0,078	300000	923800	7.868.385	82,07	83,99	80,06	300000

Bahamas, The	0,0072	18,22	8,82	2,23	16700		0,0170	0,0690		65,63	27,3	0,283	0,657	0,06	84000	131700	299.697	25,7	31,73	19,55	84000
Bahrain	0,0156	18,54	4,03	2,67	16900		-0,0020	0,1500	0,075	73,98	29	0,284	0,684	0,033	195700	185800	677.886	17,91	20,93	14,8	195700
Bangladesh	0,0208	30,03	8,52	3,15	1900	0,3560	0,0560	0,4000	0,012	61,71	21,5	0,335	0,631	0,034	243000	740000	141.340.476	64,32	65,41	63,16	243000
Barbados	0,0036	12,98	9,08	1,65	15700		-0,0050	0,1070		71,64	33,7	0,21	0,703	0,088	100000	134000	278.289	12,61	14,26	10,94	100000
Belarus	-0,0011	10,52	14,1	1,36	6100	0,2200	0,2820	0,0210	0,014	68,57	36,9	0,163	0,692	0,145	1391900	3071300	10.310.520	13,62	14,71	12,47	1391900
Belgium	0,0016	10,59	10,2	1,64	29100	0,0400	0,0160	0,0810	0,013	78,44	40,2	0,171	0,656	0,173	3400000	5120400	10.348.276	4,76	5,36	4,12	3400000
Belize	0,0239	29,89	6,04	3,77	4900	0,3300	0,0260	0,0910	0,02	67,43	19,1	0,406	0,558	0,035	30000	33300	272.945	26,37	29,75	22,82	30000
Benin	0,0289	42,57	13,69	5,95	1100	0,3700	0,0150		0,027	50,81	16,5	0,468	0,51	0,023	70000	66500	7.250.033	85,88	90,89	80,71	70000
Bermuda	0,0068	11,83	7,55	1,9	36000	0,1900	0,0330	0,0500	0,0011	77,6	39,2	0,191	0,693	0,117	34500	56000	64.935	8,79	10,45	7,1	34500
Bhutan	0,0212	34,41	13,2	4,87	1300		0,0300		0,019	53,99	20,2	0,393	0,566	0,04	15000	25200	2.185.569	102,56	100,35	104,89	15000
Bolivia	0,0156	24,65	7,77	3,08	2400	0,7000	0,0330	0,1170	0,016	65,14	21,1	0,364	0,591	0,045	270000	600100	8.724.156	54,58	58,23	50,75	270000
Bosnia and Herzegovina	0,0045	12,56	8,33	1,71	6100		0,0090	0,4000	0,045	72,57	35,9	0,189	0,706	0,105	100000	938000	4.007.608	21,88	24,5	19,08	100000
Botswana	-0,0089	24,71	33,63	3,17	9000	0,4700	0,0920	0,4000	0,036	30,76	19,2	0,392	0,562	0,046	60000	142400	1.561.973	69,98	70,96	68,96	60000
Brazil	0,0111	17,25	6,14	1,97	7600	0,2200	0,1470	0,1230	0,021	71,41	27,4	0,266	0,676	0,058	14300000	38810000	184.101.109	30,66	34,47	26,65	14300000
Brunei	0,0195	19,33	3,4	2,33	18600		-0,0200	0,1000	0,059	74,54	26,7	0,291	0,68	0,029	35000	90000	365.251	13,05	16,51	9,41	35000
Bulgaria	-0,0092	9,65	14,25	1,37	7600	0,1340	0,0230	0,1430	0,026	71,75	40,5	0,144	0,685	0,171	630000	2868200	7.517.973	21,31	25,15	17,23	630000
Burkina Faso	0,0257	44,46	18,79	6,28	1100	0,4500	0,0190		0,016	44,2	16,8	0,46	0,511	0,029	48000	65400	13.574.820	98,67	106,7	90,39	48000
Burma	0,0047	18,64	12,16	2,08	1800	0,2500	0,4970	0,0420	0,021	56,01	25,7	0,276	0,675	0,049	28000	357300	42.720.196	68,78	74,78	62,42	28000
Burundi	0,0220	39,68	17,61	5,9	600	0,6800	0,1070		0,06	43,36	16,5	0,464	0,509	0,027	14000	23900	6.231.221	70,4	77,15	63,45	14000
Cambodia	0,0180	27,13	9,1	3,51	1900	0,3600	0,0170	0,0250		58,41	19,5	0,383	0,586	0,031	30000	35400	13.363.421	73,67	82,51	64,44	30000
Cameroon	0,0197	35,08	15,34	4,55	1800	0,4800	0,0230	0,3000	0,014	47,95	18,5	0,42	0,548	0,032	60000	110900	16.063.678	69,18	73,16	65,09	60000
Canada	0,0092	10,91	7,67	1,61	29800		0,0280	0,0780	0,011	79,96	38,2	0,182	0,687	0,13	16110000	19950900	32.507.874	4,82	5,28	4,33	16110000
Central African Republic	0,0156	35,55	19,99	4,59	1100		0,0360	0,0800	0,011	41,36	18	0,428	0,538	0,034	5000	9000	3.742.482	92,15	99,09	84,99	5000
Chad	0,0300	46,5	16,38	6,38	1200	0,8000	0,0600		0,021	48,24	16	0,479	0,493	0,028	15000	11800	9.538.544	94,78	104,01	85,17	15000
Chile	0,0101	15,77	5,71	2,06	9900	0,2060	0,0280	0,0850	0,04	76,38	29,8	0,258	0,663	0,078	3575000	3467000	15.823.957	9,05	9,81	8,25	3575000

China	0,0057	12,98	6,92	1,69	5000	0,1000	0,0120	0,1010		71,96	31,8	0,223	0,703	0,075	79500000	263000000	1.298.847.624	25,28	21,84	29,14	79500000
Colombia	0,0153	21,19	5,61	2,59	6300	0,5500	0,0710	0,1420	0,034	71,43	25,8	0,31	0,639	0,05	2732200	8768100	42.310.775	21,72	25,69	17,61	2732200
Congo, Democratic Republic of the	0,0299	44,73	14,64	6,62	700		0,1400		0,014	49,14	15,8	0,482	0,493	0,025	50000	10000	58.317.930	94,69	103,18	85,95	50000
Congo, Republic of the	0,0142	28,66	14,49	3,54	700		0,0240		0,028	49,51	20,4	0,379	0,585	0,036	15000	7000	2.998.040	93,86	99,95	87,59	15000
Costa Rica	0,0152	18,99	4,32	2,33	9100	0,2060	0,0940	0,0670	0,004	76,63	25,7	0,295	0,65	0,055	800000	1132000	3.956.507	10,26	11,17	9,3	800000
Cote d'Ivoire	0,0211	39,64	18,48	5,42	1400	0,3700	0,0340	0,1300	0,012	42,48	17	0,451	0,526	0,022	90000	328000	17.327.724	97,1	113,87	79,83	90000
Croatia	-0,0002	9,51	11,3	1,39	10600		0,0180	0,1950		74,14	39,7	0,166	0,67	0,164	1014000	1825000	4.496.869	6,96	7,03	6,9	1014000
Cuba	0,0034	12,18	7,17	1,66	2900		0,0410	0,0260	0,018	77,04	34,8	0,2	0,698	0,101	120000	574400	11.308.764	6,45	7,25	5,6	120000
Czech Republic	-0,0005	9,1	10,54	1,18	15700		0,0010	0,0990	0,021	75,78	38,6	0,15	0,709	0,141	2700000	3626000	10.246.178	7,36	9,19	5,44	2700000
Denmark	0,0035	11,59	10,53	1,74	31100		0,0210	0,0610	0,016	77,44	39,2	0,189	0,662	0,15	2756000	3610100	5.413.392	3,97	4,32	3,59	2756000
Djibouti	0,0210	40,39	19,42	5,48	1300	0,5000	0,0200	0,5000	0,044	43,12	18,3	0,432	0,537	0,031	6500	9500	466.900	4,63	4,65	4,6	6500
Dominica	-0,0045	16,25	6,9	1,98	5400	0,3000	0,0100	0,2300		74,38	29	0,273	0,648	0,079	12500	23700	69.278	14,75	19,49	9,77	12500
Dominican Republic	0,0133	23,6	7,1	2,89	6000	0,2500	0,2750	0,1650	0,011	67,63	23,7	0,333	0,614	0,053	500000	901800	8.833.634	33,28	35,75	30,68	500000
East Timor	0,0211	27,46	6,36	3,7	500	0,4200	0,0800	0,5000		65,56	20	0,378	0,592	0,029	0		1.019.252	48,86	55,34	42,05	0
Ecuador	0,0103	23,18	4,26	2,78	3300	0,6500	0,0790	0,0980	0,024	76,01	23	0,339	0,612	0,049	569700	1549000	13.212.742	24,49	29,34	19,4	569700
Egypt	0,0183	23,84	5,3	2,95	4000	0,1670	0,0430	0,0990	0,036	70,71	23,4	0,334	0,622	0,043	2700000	8735700	76.117.421	33,9	34,64	33,12	2700000
El Salvador	0,0178	27,48	5,93	3,2	4800	0,4800	0,0210	0,0650	0,011	70,92	21,4	0,368	0,581	0,051	550000	752600	6.587.541	25,93	28,79	22,92	550000
Equatorial Guinea	0,0243	36,56	12,27	4,68	2700		0,0600	0,3000	0,025	55,15	18,7	0,42	0,543	0,038	1800	9600	523.051	87,08	93,27	80,71	1800
Eritrea	0,0257	39,03	13,36	5,67	700	0,5300	0,1230		0,118	52,7	17,5	0,448	0,519	0,033	9500	38100	4.447.307	75,59	83,03	67,94	9500
Estonia	-0,0066	9,79	13,27	1,39	12300		0,0130	0,1010		71,38	38,8	0,16	0,675	0,165	444000	475000	1.341.664	8,08	9,32	6,76	444000
Ethiopia	0,0189	39,23	20,36	5,44	700	0,5000	0,1780		0,052	40,88	17,4	0,447	0,525	0,028	75000	435000	67.851.281	102,12	112,22	91,72	75000
Faroe Islands	0,0066	13,89	8,68	2,22	22000		0,0510	0,0100		79,05	35,1	0,216	0,646	0,137	25000	23000	46.662	6,38	7,72	5,03	25000
Fiji	0,0141	22,91	5,68	2,78	5800	0,2550	0,0160	0,0760	0,022	69,2	24	0,317	0,643	0,04	55000	102000	880.874	12,99	14,36	11,54	55000

Finland	0,0018	10,56	9,69	1,73	27400		0,0090	0,0900	0,02	78,24	40,7	0,175	0,667	0,157	2650000	2548000	5.214.512	3,59	3,91	3,26	2650000
France	0,0039	12,34	9,06	1,85	27600	0,0650	0,0210	0,0970	0,026	79,44	38,6	0,185	0,651	0,164	21900000	33905400	60.424.213	4,31	4,83	3,78	21900000
French Guiana	0,0225	21	4,82	3,05	8300		0,0150	0,2200		76,89	28,3	0,296	0,645	0,059	3200	51000	191.309	12,46	13,3	11,58	3200
Gabon	0,0250	36,4	11,43	4,8	5500		0,0050	0,2100	0,02	56,46	18,5	0,422	0,537	0,041	35000	38400	1.355.246	54,34	64,15	44,23	35000
Gambia, The	0,0298	40,3	12,08	5,46	1700		0,1400		0,003	54,79	17,5	0,447	0,526	0,027	25000	38400	1.546.848	73,48	80,14	66,62	25000
Gaza Strip	0,0383	40,62	3,95	6,04	600	0,6000	0,0220			71,59	15,5	0,49	0,483	0,027	60000	95729	1.324.991	23,54	24,71	22,31	60000
Georgia	-0,0036	10,1	8,98	1,4	2500	0,5400	0,0480	0,1700	0,0059	75,62	37	0,187	0,658	0,155	150500	650500	4.693.892	19,34	21,5	16,84	150500
Germany	0,0002	8,45	10,44	1,38	27600		0,0110	0,1050	0,015	78,54	41,7	0,147	0,67	0,183	39000000	54350000	82.424.609	4,2	4,64	3,73	39000000
Ghana	0,0136	24,9	10,67	3,17	2200	0,3140	0,2670	0,2000	0,006	56,27	20	0,38	0,583	0,037	170000	302300	20.757.032	52,22	55,1	49,25	170000
Greece	0,0020	9,73	10,08	1,32	20000		0,0360	0,0940	0,043	78,94	40,2	0,145	0,67	0,186	1718400	5205100	10.647.529	5,63	6,19	5,04	1718400
Grenada	0,0014	22,61	7,31	2,41	5000	0,3200	0,0280	0,1250		64,52	20,9	0,345	0,62	0,035	15000	33500	89.357	14,62	14,18	15,07	15000
Guadeloupe	0,0096	15,79	6,05	1,91	8000			0,2780		77,71	31,4	0,244	0,666	0,09	20000	210000	444.515	8,83	10,07	7,53	20000
Guam	0,0150	19,31	4,35	2,61	21000	0,2300	0,0000	0,1500		78,12	28,2	0,298	0,64	0,063	50000	84134	166.090	7,15	7,8	6,46	50000
Guatemala	0,0261	34,58	6,79	4,6	4100	0,7500	0,0550	0,0750	0,008	65,19	18,4	0,426	0,54	0,033	400000	846000	14.280.596	36,91	37,71	36,07	400000
Guinea	0,0237	42,26	15,53	5,87	2100	0,4000	0,1480		0,017	49,7	17,7	0,444	0,524	0,032	40000	26200	9.246.462	91,82	97,3	86,18	40000
Guinea-Bissau	0,0199	38,03	16,57	5	800		0,0400		0,028	46,98	18,9	0,417	0,554	0,029	19000	10600	1.388.363	108,72	119,37	97,74	19000
Guyana	0,0061	17,85	9,71	2,06	4000		0,0570	0,0910	0,008	62,43	26,2	0,265	0,683	0,051	125000	80400	705.803	37,22	41,28	32,95	125000
Haiti	0,0171	33,76	13,21	4,76	1600	0,8000	0,3780		0,009	51,78	18,1	0,422	0,541	0,037	80000	130000	7.656.166	74,38	79,83	68,65	80000
Honduras	0,0224	31,04	6,64	3,97	2600	0,5300	0,0770	0,2750	0,015	66,15	19	0,412	0,551	0,037	168600	322500	6.823.568	29,64	33,22	25,89	168600
Hong Kong	0,0065	7,23	5,98	0,91	28800		-0,0260	0,0790		81,39	39,4	0,142	0,733	0,125	3212800	3801300	6.855.125	2,97	3,16	2,77	3212800
Hungary	-0,0025	9,77	13,16	1,31	13900	0,0860	0,0470	0,0590		72,25	38,4	0,16	0,69	0,15	1600000	3666400	10.032.375	8,68	9,34	7,98	1600000
Iceland	0,0097	13,83	6,57	1,93	30900		0,0210	0,0340		80,18	33,8	0,225	0,658	0,117	195000	190700	293.966	3,31	3,44	3,17	195000
India	0,0144	22,8	8,38	2,85	2900	0,2500	0,0380	0,0950	0,024	63,99	24,4	0,317	0,635	0,048	18481000	48917000	1.065.070.607	57,92	58,52	57,29	18481000
Indonesia	0,0149	21,11	6,26	2,47	3200	0,2700	0,0660	0,0870	0,013	69,26	26,1	0,294	0,655	0,051	80000	7750000	238.452.952	36,82	42,09	31,29	80000
Iran	0,0107	17,1	5,53	1,93	7000	0,4000	0,1640	0,1570		69,66	23,5	0,28	0,672	0,048	4300000	14571100	69.018.924	42,86	43,01	42,69	4300000
Iraq	0,0274	33,09	5,66	4,4	1500		0,2930			68,26	19,2	0,403	0,567	0,03	25000	675000	25.374.691	52,71	58,58	46,55	25000

Ireland	0,0116	14,47	7,91	1,87	29600	0,1000	0,0350	0,0470	0,009	77,36	33,4	0,21	0,675	0,115	1260000	1955000	3.969.558	5,5	6,04	4,92	1260000
Israel	0,0129	18,45	6,19	2,47	19800	0,1800	0,0070	0,1070	0,087	79,17	29,2	0,267	0,634	0,099	20000	3006000	6.199.008	7,21	7,96	6,42	20000
Italy	0,0009	9,05	10,21	1,27	26700		0,0270	0,0860	0,019	79,54	41,4	0,14	0,669	0,191	18500000	26596000	58.057.477	6,07	6,68	5,41	18500000
Jamaica	0,0066	16,94	5,4	1,98	3900	0,1970	0,1030	0,1590	0,004	76,07	26,8	0,282	0,65	0,069	600000	444400	2.713.130	12,81	13,82	11,75	600000
Japan	0,0008	9,56	8,75	1,38	28200		-0,0030	0,0530	0,01	81,04	42,3	0,143	0,667	0,19	57200000	71149000	127.333.002	3,28	3,54	3	57200000
Jordan	0,0267	22,73	2,62	2,86	4300	0,3000	0,0240	0,2500	0,202	78,06	22,2	0,352	0,611	0,037	457000	622600	5.611.202	18,11	21,63	14,37	457000
Kazakhstan	0,0026	15,52	9,59	1,9	6300	0,2600	0,0660	0,0880		66,07	28,3	0,244	0,68	0,076	250000	2081900	15.143.704	30,54	35,24	25,57	250000
Kenya	0,0114	27,82	16,31	3,31	1000	0,5000	0,0980	0,4000	0,018	44,94	18,6	0,406	0,565	0,029	400000	328400	32.021.856	62,62	65,55	59,6	400000
Korea, North	0,0098	16,77	6,99	2,2	1300				0,229	71,08	31,4	0,246	0,678	0,076		1100000	22.697.553	24,84	26,59	23	0
Korea, South	0,0062	12,33	6,13	1,56	17800	0,0400	0,0360	0,0340	0,027	75,58	33,7	0,204	0,714	0,082	29220000	22877000	48.598.175	7,18	7,64	6,68	29220000
Kuwait	0,0336	21,85	2,44	3,03	19000		0,0120	0,0210	0,058	76,84	25,9	0,275	0,698	0,027	567000	486900	2.257.549	10,26	11,27	9,22	567000
Kyrgyzstan	0,0125	22,13	7,19	2,71	1600	0,5000	0,0350	0,0720	0,014	67,84	23,1	0,323	0,616	0,061	152000	394800	5.081.429	36,81	42,23	31,11	152000
Laos	0,0244	36,47	12,1	4,86	1700	0,4000	0,1530	0,0570	0,005	54,69	18,6	0,419	0,549	0,032	15000	61900	6.068.117	87,06	97,05	76,69	15000
Latvia	-0,0071	8,87	13,73	1,25	10200		0,0290	0,0860	0,012	70,86	38,8	0,15	0,692	0,158	936000	653900	2.306.306	9,67	11,45	7,8	936000
Lebanon	0,0130	19,31	6,28	1,95	4800	0,2800	0,0250	0,1800	0,048	72,35	26,9	0,269	0,663	0,069	400000	678800	3.777.218	25,48	28,21	22,61	400000
Lesotho	0,0014	26,91	24,79	3,44	3000	0,4900	0,0610	0,4500	0,026	36,81	20	0,373	0,572	0,055	21000	28600	1.865.040	85,22	90,19	80,1	21000
Liberia	0,0270	44,81	17,86	6,16	1000	0,8000	0,1500	0,8500	0,013	47,93	18,1	0,434	0,529	0,036	1000	7000	3.390.635	130,51	137,32	123,5	1000
Libya	0,0237	27,17	3,48	3,42	6400		0,0280	0,3000	0,039	76,28	22,4	0,342	0,617	0,041	160000	750000	5.631.585	25,7	28,04	23,26	160000
Liechtenstein	0,0086	10,65	6,94	1,51	25000		0,0100	0,0130		79,4	38,8	0,178	0,705	0,116	20000	19900	33.436	4,77	6,47	3,07	20000
Lithuania	-0,0033	8,49	11,03	1,17	11400		-0,0120	0,1030	0,019	73,46	37,4	0,167	0,684	0,149	695700	824200	3.607.899	7,13	8,61	5,57	695700
Luxembourg	0,0128	12,21	8,42	1,79	55100		0,0200	0,0360	0,009	78,58	38,3	0,19	0,664	0,145	165000	355400	462.690	4,88	4,84	4,92	165000
Macau	0,0087	8,04	4,23	0,93	19400		-0,0260	0,0630		82,03	35,2	0,18	0,742	0,078	120000	174600	445.286	4,39	4,59	4,19	120000
Macedonia, The Former Yugoslav Republic of	0,0039	13,14	7,83	1,74	6700	0,3020	0,0120	0,3670		74,73	32,8	0,215	0,678	0,107	100000	560000	2.071.210	11,74	12,67	10,73	100000
Madagascar	0,0303	41,91	11,62	5,7	800	0,7100	0,0800	0,0590	0,012	56,54	17,4	0,449	0,52	0,031	70500	59600	17.501.871	78,52	86,84	69,95	70500
Malawi	0,0214	44,35	23,01	6,04	600	0,5500	0,0950		0,007	37,48	16,4	0,468	0,505	0,028	36000	85000	11.906.855	104,23	108,4	99,93	36000

Malaysia	0,0183	23,37	5,08	3,1	9000	0,0800	0,0110	0,0360	0,0203	71,95	23,8	0,333	0,621	0,045	8692100	4571600	23.522.482	18,35	21,23	15,27	8692100
Maldives	0,0286	36,06	7,44	5,14	3900		0,0100		0,086	63,68	17,5	0,444	0,526	0,03	15000	28700	339.330	58,32	57,43	59,26	15000
Mali	0,0278	47,29	19,12	6,58	900	0,6400	0,0450	0,1460	0,013	45,28	16,3	0,471	0,499	0,03	25000	56600	11.956.788	117,99	124,51	111,27	25000
Malta	0,0042	10,09	7,93	1,49	17700		0,0040	0,0700	0,007	78,68	38	0,181	0,685	0,134	120000	208300	396.851	3,94	4,45	3,4	120000
Martinique	0,0081	14,56	6,43	1,79	14400		0,0390	0,2720		78,88	33,1	0,226	0,67	0,104	40000	172000	429.510	7,27	4,79	9,8	40000
Mauritania	0,0291	41,79	12,74	6,01	1800	0,5000	0,0700	0,2100	0,037	52,32	16,9	0,459	0,519	0,022	10000	31500	2.998.563	72,35	75,22	69,39	10000
Mauritius	0,0081	15,85	6,82	1,97	11400	0,1000	0,0420	0,0980	0,002	72,09	30,1	0,248	0,688	0,065	150000	348200	1.220.481	15,57	18,36	12,73	150000
Mexico	0,0118	21,44	4,73	2,49	9000	0,4000	0,0450	0,0330	0,009	74,94	24,6	0,316	0,629	0,055	10033000	15958700	104.959.594	21,69	23,63	19,65	10033000
Micronesia, Federated States of	-0,0002	25,8	4,97	3,37	2000	0,2670	0,0100	0,1600		69,44		0,377	0,592	0,031	6000	10100	108.155	31,28	34,45	27,95	6000
Moldova	0,0018	14,81	12,76	1,78	1800	0,8000	0,1160	0,0800	0,004	65,03	32,1	0,206	0,691	0,103	150000	706900	4.446.455	41	43,96	37,89	150000
Monaco	0,0044	9,36	12,74	1,76	27000			0,0310		79,42	45,1	0,155	0,621	0,224	16000	33700	32.270	5,53	6,4	4,61	16000
Mongolia	0,0143	21,44	7,1	2,27	1800	0,3300	0,0150	0,0460	0,022	64,17	23,9	0,297	0,667	0,036	50000	128000	2.751.314	55,45	58,97	51,75	50000
Morocco	0,0161	22,79	5,71	2,81	4000	0,1900	0,0120	0,1900	0,048	70,35	23,3	0,326	0,625	0,049	800000	1219200	32.209.101	43,25	47,09	39,21	800000
Mozambique	0,0122	36,06	23,86	4,78	1200	0,7000	0,1400	0,2100	0,022	37,1	18,2	0,436	0,536	0,028	50000	83700	18.811.731	137,08	142,67	131,32	50000
Namibia	0,0125	33,51	21,02	4,65	7200	0,5000	0,0730	0,3500	0,025	40,53	18,3	0,424	0,54	0,035	65000	127400	1.954.033	69,58	72,65	66,43	65000
Nepal	0,0223	31,96	9,66	4,29	1400	0,4200	0,0290	0,4700	0,016	59,4	19,9	0,394	0,57	0,036	80000	371800	27.070.666	68,77	67,1	70,53	80000
Netherlands	0,0057	11,41	8,67	1,66	28600		0,0210	0,0530	0,016	78,68	38,7	0,183	0,678	0,139	8500000	10004000	16.318.199	5,11	5,68	4,51	8500000
New Caledonia	0,0133	18,98	5,64	2,35	15000		-0,0060	0,1900		73,78	27,2	0,294	0,643	0,063	60000	52000	213.679	7,89	8,59	7,14	60000
New Zealand	0,0105	14,04	7,54	1,79	21600		0,0180	0,0470	0,01	78,49	33,4	0,217	0,667	0,116	2110000	1765000	3.993.817	5,96	6,83	5,05	2110000
Nicaragua	0,0197	25,5	4,54	2,89	2300	0,5000	0,0530	0,2200	0,012	70,02	20,2	0,381	0,589	0,03	90000	171600	5.359.759	30,15	33,73	26,38	90000
Niger	0,0267	48,91	21,51	6,83	800	0,6300	0,0300		0,011	42,18	16,2	0,475	0,504	0,021	15000	22400	11.360.538	122,66	126,96	118,24	15000
Nigeria	0,0245	38,24	13,99	5,32	900	0,6000	0,1380	0,0270	0,009	50,49	18,1	0,434	0,537	0,029	750000	853100	137.253.133	70,49	73,55	67,34	750000
Northern Mariana Islands	0,0271	19,77	2,3	1,33	12500		0,0120			75,67	29,2	0,203	0,781	0,016		21000	78.252	7,25	7,18	7,32	0
Norway	0,0041	11,89	9,51	1,78	37800		0,0250	0,0470	0,019	79,25	37,9	0,198	0,654	0,148	2288000	3343000	4.574.560	3,73	4,11	3,34	2288000

Oman	0,0335	37,12	3,91	5,9	13100		-0,0030		0,114	72,85	19,3	0,424	0,551	0,025	180000	233900	2.903.165	20,26	23,18	17,19	180000
Pakistan	0,0198	31,22	8,67	4,29	2100	0,3500	0,0290	0,0770	0,039	62,61	19,4	0,402	0,558	0,041	1500000	3982800	159.196.336	74,43	74,84	74	1500000
Panama	0,0131	20,36	6,39	2,49	6300	0,3700	0,0140	0,1380	0,012	72,14	25,9	0,302	0,636	0,062	120000	386900	3.000.463	20,95	23,08	18,72	120000
Papua New Guinea	0,0230	30,52	7,5	4,04	2200	0,3700	0,1470		0,014	64,56	21	0,383	0,58	0,038	75000	62000	5.420.280	53,15	57,39	48,7	75000
Paraguay	0,0251	29,78	4,58	3,98	4700	0,3600	0,1420	0,1850	0,009	74,64	21,1	0,382	0,57	0,048	120000	273200	6.191.368	26,67	31,5	21,6	120000
Peru	0,0139	21,27	6,29	2,61	5100	0,5400	0,0230	0,0970	0,013	69,22	24,6	0,321	0,628	0,051	2850000	1839200	27.544.305	32,95	35,57	30,21	2850000
Philippines	0,0188	25,8	5,53	3,22	4600	0,4000	0,0310	0,1140	0,015	69,6	22,1	0,358	0,602	0,039	3500000	3310900	86.241.697	24,24	27,11	21,22	3500000
Poland	0,0002	10,64	9,97	1,38	11100	0,1840	0,0070	0,2000	0,0171	74,16	36,2	0,171	0,7	0,129	8970000	12300000	38.626.349	8,73	9,82	7,59	8970000
Portugal	0,0041	10,9	10,37	1,46	18000		0,0330	0,0640	0,023	77,35	37,9	0,167	0,664	0,169	3600000	4278800	10.524.145	5,13	5,6	4,62	3600000
Puerto Rico	0,0049	14,1	7,77	1,91	16800		0,0650	0,1200		77,49	33,8	0,224	0,654	0,122	600000	1329500	3.897.960	8,37	9,7	6,98	600000
Qatar	0,0274	15,6	4,52	2,95	21500		0,0230	0,0270	0,1	73,4	31,4	0,242	0,727	0,031	126000	184500	840.290	19,32	22,77	15,7	126000
Romania	-0,0011	10,69	11,69	1,35	7000	0,4450	0,1530	0,0720	0,0247	71,12	36,1	0,162	0,694	0,144	4000000	4300000	22.355.551	27,24	30,41	23,86	4000000
Russia	-0,0045	9,63	15,17	1,26	8900	0,2500	0,1370	0,0850		66,39	37,9	0,15	0,713	0,137	6000000	35500000	143.782.338	16,96	19,58	14,18	6000000
Rwanda	0,0182	40,01	21,86	5,55	1300	0,6000	0,0750		0,029	39,18	18,2	0,423	0,55	0,027	25000	23200	7.954.013	101,68	106,68	96,53	25000
Samoa	-0,0025	15,69	6,47	3,11	5600		0,0400			70,41	24,2	0,283	0,655	0,063	4000	11800	177.714	28,72	33,83	23,35	4000
San Marino	0,0133	10,31	7,96	1,32	34600		0,0330	0,0260		81,53	40	0,165	0,668	0,167	14300	20600	28.503	5,85	6,28	5,38	14300
Saudi Arabia	0,0244	29,74	2,66	4,11	11800		0,0050	0,2500	0,1	75,23	21,2	0,383	0,593	0,023	1500000	3502600	25.795.938	13,7	15,72	11,58	1500000
Senegal	0,0252	35,72	10,74	4,84	1600	0,5400	0,0000	0,4800	0,015	56,56	18	0,432	0,537	0,03	225000	228800	10.852.147	56,53	60,25	52,71	225000
Serbia and Montenegro	0,0003	12,13	10,53	1,67	2200	0,3000	0,1120	0,3450		74,4	36,6	0,183	0,668	0,149	847000	2611700	10.825.900	13,43	15,04	11,68	847000
Seychelles	0,0045	16,55	6,41	1,77	7800		0,0330		0,018	71,53	27,3	0,269	0,669	0,062	11700	21700	80.832	15,97	20,2	11,61	11700
Sierra Leone	0,0227	43,34	20,62	5,79	500	0,6800	0,0100		0,015	42,69	17,5	0,448	0,52	0,032	8000	24000	5.883.889	145,24	162,64	127,32	8000
Singapore	0,0171	9,63	4,05	1,04	23700		0,0050	0,0480	0,049	81,53	36,2	0,165	0,758	0,078	2310000	1896100	4.353.893	2,28	2,49	2,05	2310000
Slovakia	0,0014	10,57	9,48	1,31	13300		0,0860	0,1520	0,0189	74,19	35,1	0,175	0,708	0,117	1375800	1294700	5.423.567	7,62	8,88	6,28	1375800
Slovenia	-0,0001	8,9	10,15	1,23	19000		0,0560	0,1120	0,017	75,93	39,8	0,143	0,706	0,151	750000	812300	2.011.473	4,5	5,11	3,84	750000
Somalia	0,0341	46,04	17,3	6,91	500				0,009	47,71	17,6	0,447	0,527	0,027	89000	100000	8.304.601	118,52	127,95	108,81	89000
South Africa	-0,0025	18,38	20,54	2,18	10700	0,5000	0,0590	0,3100	0,017	44,19	24,7	0,295	0,653	0,052	3100000	4844000	42.718.530	62,18	65,87	58,4	3100000

Spain	0,0016	10,11	9,55	1,27	22000		0,0300	0,1130	0,012	79,37	39,1	0,144	0,68	0,176	9789000	17567500	40.280.780	4,48	4,88	4,06	9789000
Sri Lanka	0,0081	15,88	6,47	1,88	3700	0,2200	0,0630	0,0840	0,032	72,89	29,1	0,248	0,682	0,07	200000	881400	19.905.165	14,78	16,01	13,5	200000
Sudan	0,0264	35,79	9,37	4,97	1900		0,0880	0,1870	0,025	58,13	17,9	0,437	0,541	0,023	300000	900000	39.148.162	64,05	64,8	63,26	300000
Suriname	0,0031	18,87	6,99	2,37	4000	0,7000	0,1700	0,1700	0,007	69,1	25,8	0,302	0,637	0,061	20000	79800	436.935	24,15	28,24	19,85	20000
Swaziland	0,0055	28,55	23,06	3,81	4900	0,4000	0,0730	0,3400	0,018	37,54	18,6	0,41	0,553	0,037	27000	46200	1.169.241	68,35	71,64	64,96	27000
Sweden	0,0018	10,46	10,38	1,66	26800		0,0190	0,0490	0,021	80,3	40,3	0,175	0,652	0,173	5125000	6579200	8.986.400	2,77	2,93	2,6	5125000
Switzerland	0,0054	9,83	8,44	1,42	32700		0,0060	0,0370	0,01	80,31	39,5	0,168	0,679	0,153	2556000	5419000	7.450.867	4,43	4,94	3,89	2556000
Syria	0,0240	28,93	4,96	3,61	3300	0,2000	0,0150	0,2000	0,059	69,71	20	0,38	0,587	0,033	220000	2099300	18.016.874	30,6	30,82	30,35	220000
Taiwan	0,0064	12,7	6,29	1,57	23400	0,0100	-0,0030	0,0500	0,027	77,06	33,7	0,199	0,707	0,094	8830000	13355000	22.749.838	6,52	7,21	5,77	8830000
Tajikistan	0,0214	32,63	8,42	4,11	1000	0,6000	0,1630	0,4000	0,039	64,47	19,5	0,392	0,561	0,047	4100	242100	7.011.556	112,1	124,47	99,11	4100
Tanzania	0,0195	39	17,45	5,15	600	0,3600	0,0440		0,002	44,39	17,6	0,442	0,532	0,026	250000	149100	36.588.225	102,13	111,62	92,35	250000
Thailand	0,0091	16,04	6,94	1,89	7400	0,1040	0,0180	0,0220	0,018	71,41	30,5	0,241	0,687	0,073	6031300	6600000	64.865.523	21,14	22,49	19,73	6031300
Togo	0,0227	34,36	11,64	4,79	1500	0,3200	-0,0100		0,019	53,05	17,6	0,439	0,536	0,026	210000	60600	5.556.812	67,66	75,4	59,69	210000
Tonga	0,0194	24,87	5,45	3	2200		0,1030	0,1330		69,2	20,1	0,371	0,587	0,042	2900	11200	110.237	12,99	14,36	11,54	2900
Trinidad and Tobago	-0,0071	12,75	9,02	1,77	9500	0,2100	0,0380	0,1040	0,006	69,28	30,4	0,214	0,706	0,081	138000	325100	1.096.585	24,64	26,58	22,61	138000
Tunisia	0,0101	15,74	5,05	1,79	6900	0,0760	0,0270	0,1430	0,015	74,66	26,8	0,26	0,675	0,065	630000	1163800	9.974.722	25,76	28,7	22,59	630000
Turkey	0,0113	17,22	5,95	1,98	6700	0,1800	0,2530	0,1050	0,053	72,08	27,3	0,266	0,668	0,066	5500000	18916700	68.893.918	42,62	46,3	38,76	5500000
Turkmenistan	0,0181	27,82	8,82	3,45	5800	0,3440	0,0950		0,034	61,29	21,3	0,362	0,597	0,041	8000	374000	4.863.169	73,13	76,9	69,16	8000
Uganda	0,0297	46,31	16,61	6,64	1400	0,3500	0,0790		0,021	45,28	14,8	0,506	0,471	0,024	125000	61000	26.404.543	86,15	93,58	78,5	125000
Ukraine	-0,0066	10,21	16,41	1,37	5400	0,2900	0,0520	0,0370	0,014	66,68	38,1	0,159	0,687	0,154	900000	10833300	47.732.079	20,61	21,87	19,28	900000
United Arab Emirates	0,0157	18,65	4,14	3,02	23200		0,0320	0,0240	0,031	74,99	27,7	0,259	0,709	0,032	1110200	1135800	2.523.915	15,06	17,71	12,29	1110200
United Kingdom	0,0029	10,88	10,19	1,66	27700	0,1700	0,0140	0,0500	0,024	78,27	38,7	0,18	0,663	0,157	250000	34898000	60.270.708	5,22	5,83	4,58	250000
United States	0,0092	14,13	8,34	2,07	37800	0,1200	0,0230	0,0600		77,43	36	0,208	0,669	0,124	159000000	181599900	293.027.571	6,63	7,31	5,91	159000000
Uruguay	0,0051	14,44	9,07	1,96	12800	0,2370	0,1940	0,1600	0,02	75,92	32,2	0,235	0,634	0,131	400000	946500	3.399.237	12,31	13,67	10,91	400000
Uzbekistan	0,0165	26,12	7,95	2,97	1700		0,1310	0,0050	0,02	64,09	22,1	0,341	0,611	0,048	492000	1717100	26.410.416	71,3	75,03	67,39	492000

Venezuela	0,0144	19,34	4,9	2,31	4800	0,4700	0,3110	0,1800	0,013	74,06	25,2	0,305	0,645	0,05	1274400	2841800	25.017.387	22,99	26,18	19,57	1274400
Vietnam	0,0130	19,58	6,14	2,22	2500	0,3700	0,0310	0,0610	0,025	70,35	24,9	0,294	0,65	0,056	3500000	4402000	82.689.518	29,88	33,71	25,77	3500000
West Bank	0,0321	33,21	4,07	4,52	800	0,6000	0,0220	0,5000		72,88	18	0,438	0,528	0,035	145000	301600	2.311.204	20,16	22,28	17,91	145000
Yemen	0,0344	43,16	8,78	6,75	800	0,1570	0,1080	0,3500	0,079	61,36	16,5	0,466	0,506	0,028	100000	542200	20.024.867	63,26	68,12	58,15	100000
Zambia	0,0147	38,99	24,35	5,14	800	0,8600	0,2140	0,5000	0,009	35,18	16,6	0,461	0,511	0,028	68200	88400	10.462.436	98,4	105,6	90,98	68200
Zimbabwe	0,0068	30,05	23,3	3,6	1900	0,7000	3,8470	0,7000	0,017	37,82	19,1	0,394	0,57	0,036	500000	300900	12.671.860	67,08	69,74	64,33	500000

12.

Sandkasten für Übungsaufgaben

[<< hoch zum Anfang vom Inhaltsverzeichnis](#)

[<< hoch zum Kapitel 12 vom Inhaltsverzeichnis](#)

Sandkasten für Aufgabenvorschläge, Kritik, Lösungen usw.

Das Ganze entwickelt eine unvorhergesehene Eigendynamik, was zwar chaotisch, aber auch durchaus reizvoll ist. Mal schauen, wann die Diskussionsseite in Betrieb genommen werden muss

Bemerkungen

Eigentlich dachte ich hier an neue Aufgaben ... Aber trotzdem danke, der Hinweis in Aufgabe 2 ist richtig.

Hallo Unbekannte(r), da Du ja offensichtlich gern Aufgaben löst, schau doch mal weitere Aufgaben in [Mathematik: Statistik: Übungsaufgaben](#) kritisch durch. Ich habe übrigens die bereits von Dir gelösten noch nicht überprüft, weil ich dann a) die Lösungen hervorkramen müsste oder b) selber nachrechnen müsste, wozu bei beidem ich noch nicht gekommen bin. Scheinen aber auf den ersten Blick o.k. zu sein. --[Philipendula](#) 12:30, 12. Sep 2004 (UTC)

Also Aufgabe 1 ist falsch ... $\Omega = \{1,2,3,4,5,6,7\}$ - dabei bezeichnet $\{1\} =$ "die Anlage wird in einem Tag fertig" etc.; $A = \{5,6,7\}$ $B = \{1,2,3,4,5\}$ 1. $A_c = \{1,2,3,4\}$ d.h. die Anlage wird in höchstens 4 Tagen fertig 2. Die Schnittmenge zweier Mengen bezeichnet die Elemente, die beide Mengen gemeinsamen haben: $A \cap B = \{5\}$ 3. Ist diese Schnittmenge leer, heißen die Mengen disjunkt (kleine Ergänzung:-): also nein, wie bereits angegeben. 4. sowohl $A \cap B$ als auch A_c sind Teilmengen von B . bzw: $\{5\}$ vereinigt $\{1,2,3,4\}$ vereinigt $B = B$ vereinigt $B = B$.

Bei Aufgabe 5 ist $\max(P(L|V), P(S|V), P(F|V))$ gesucht, da die Vorbedingung lautet, dass die Wand zerkratzt wurde. Wenn jeder nur ~2% verdächtig wäre, wäre das sehr mystisch.

Jo. Wurde von einer IP gelöst. Ich denke, die Person hat sich auf den Zähler der W. konzentriert, denn es langt ja, wenn man den maximalen Zähler ermittelt. --[Philipendula](#) 16:06, 27. Apr 2006 (UTC)

Platz für Aufgabenvorschläge

Hier war schon jemand mutig:

Aufgabe 1

Eine Porzellanfabrik erhält eine neue Mischanlage für spezielles Steingut. Diese muss eingerichtet und angepasst werden. Man geht davon aus, dass die Anlage in höchstens neun Tagen einsatzbereit ist. Wir definieren als Ereignisse

A: Es dauert mehr als 4 Tage, bis die Anlage einsatzbereit ist. B: Es dauert weniger als 6 Tage, bis die Anlage einsatzbereit ist.

1. Beschreiben Sie das Komplement zu A.
2. Beschreiben Sie die Schnittmenge zwischen A und B.
3. Sind A und B disjunkt?
4. Zeigen Sie, dass $(A \cap B) \cup (\text{nicht } A \cap B) = B$ ist.

Aufgabe 2

Wir beziehen uns auf Aufgabe 1 Die Werksleitung vermutet für die Zahl der Tage, die benötigt werden, um die Anlage einzurichten, die Wahrscheinlichkeiten, wie in der folgenden Tabelle angegeben:

Zahl der Tage	3	4	5	6	7
Wahrscheinlichkeit	0,08	0,24	0,41	0,20	0,07

1. Geben Sie die Wahrscheinlichkeiten für A und B an.
2. Geben Sie die Wahrscheinlichkeit für die Schnittmenge von A und B an.
3. Geben Sie die Wahrscheinlichkeit für die Vereinigungsmenge von A und B an.
4. Jeder unproduktive Tag kostet die Firma 3000 Euro. Mit welcher Wahrscheinlichkeit muss die Firma mit Kosten von höchstens 15.000 Euro rechnen?

Aufgabe 3

Die Leitung eines Kaufhauses hat für das vergangene Jahr die Zahl der wöchentlich eingegangenen Beschwerden über Servicemängel im Kaufhaus analysiert. Es ergaben sich die Wahrscheinlichkeiten für die Zahl der Beschwerden pro Woche:

Zahl der Beschwerden	0	1 - 3	4 - 6	7 - 9	10 - 12	mehr als 12
Wahrscheinlichkeit	0,14	0,39	0,23	0,15	0,06	0,03

Wir definieren die Ereignisse:

- A: Es trifft in einer Woche mindestens eine Beschwerde ein
 B: Es treffen in einer Woche weniger als 10 Beschwerden ein.

1. Geben Sie die Wahrscheinlichkeiten für A und B an.
2. Bestimmen Sie die Gegenwahrscheinlichkeit von A.
3. Geben Sie die Wahrscheinlichkeit für die Schnittmenge von A und B an.
4. Geben Sie die Wahrscheinlichkeit für die Vereinigungsmenge von A und B an.

Aufgabe 4

Die Auto-Alarmanlage Heulomat heult erfahrungsgemäß bei 90% der Autoknacker, die sich am Auto zu schaffen machen. Leider heult sie auch bei 60% aller harmlosen Kollisionen, beispielsweise mit Spaziergängern. Man vermutet, dass insgesamt 80% aller Erschütterungen eines Autos harmlos sind.

1. In wie viel Prozent aller Fälle heult die Anlage berechtigterweise?
2. Wie groß ist die Wahrscheinlichkeit, dass bei einer beliebigen Erschütterung die Anlage nicht heult?

Aufgabe 5

In einem Mietshaus wird Dienstags die Mülltonne entleert. Bei 30% der Leerungen stellt Herr Löhlein die Mülltonne raus, bei 20% der Leerungen Frau Susemihl und bei 50% aller Leerungen Herr Feinbein. Eines Tages stellt der Hausmeister fest, dass die Wand im Flur verschrämmt ist. Er weiß, dass Herr Löhlein beim Mülltonne Tragen mit einer Wahrscheinlichkeit von 7%, Frau Susemihl mit einer Wahrscheinlichkeit von 8% und Herr Feinbein mit einer Wahrscheinlichkeit von 5% mit der Tonne an der Wand entlang kratzen.

1. Welcher Bewohner ist am „verdächtigsten“?
2. Mit welcher Wahrscheinlichkeit wird nächsten Dienstag die Wand verkratzt?

Aufgabe 6

Ein Zulieferer von mechanischen Teilen für die Autoindustrie hat sich in den Liefervereinbarungen verpflichtet, die vereinbarten Qualitätsstandards regelmässig zu kontrollieren und mit einer Sicherheit von mindestens 95% einzuhalten. Für den Durchmesser des Wellentyps W1 gilt als Qualitätsstandard, dass er normal verteilt sein soll mit einem Mittelwert von 200 mm und einer Standardabweichung von 5 mm. Dabei sind Abweichungen in beiden Richtungen überprüfungsrelevant und sollen mit einer Genauigkeit von 1/10 mm berücksichtigt werden. Die Qualitätskontrolle findet standardmässig mit Zufallstichproben von 100 Wellen statt. Im vorliegenden Fall soll die Einhaltung des vereinbarten Mittelwertes überprüft werden.

1. Formulieren sie die Hypothesen
2. nennen sie die Prüfgröße und ihre Wahrscheinlichkeitsverteilung (verteilungstyp, Parameter und -werte)
3. stellen sie die prüfverteilung in geeigneter form graphisch dar
4. ermitteln sie nachvollziehbar den annahme- und ablehnbereich des tests.

Lösungen

Aufgabe 1

1. = (nicht A) = *höchstens 6 Tage* → "mehr als 6" : 7,8,9 !
2. = (A und B) = 7 Tage
3. *Nein*, sie haben *ein gemeinsames Element*.
4. (A und B) = 7d; (nichtA und B) = ($\leq 6d$ und $< 8d$) = (nicht A) = ($\leq 6d$); \Rightarrow (7d und $\leq 6d$) = $< 8d$ = B !!!

Aufgabe 2

1. $w(A) = w(5) + w(6) + w(7) = 0,68$; $w(B) = w(3) + w(4) + w(5) = 0,73$; $w(A \text{ und } B) = w(5) = 0,41$ (Aufgabe unpräzise formuliert - Wahrscheinlichkeiten für A **und** für B oder W.-keit für (A und B) ?)
2. $w(A \text{ und } B) = w(5) = 0,41$
3. $w(A \text{ oder } B) = w(A) + w(B) - w(A \text{ und } B) = 1$
4. $w(k \leq 15000) = w(k \leq (15000:3000)) = w(k \leq 5) = w(B) = 0,73$

Aufgabe 3

1. $w(A) = 1 - w(0) = 0,86$; $w(B) = 1 - (w(10-12) + w(>12)) = 0,91$
2. $w(\text{nicht}A) = 1 - w(A) = w(0) = 0,14$
3. $w(A \text{ und } B) = 1 - (w(\text{nicht}A) + w(\text{nicht}B)) = 0,77$
4. $w(A \text{ oder } B) = w(A) + w(B) - w(A \text{ und } B) = 1$

geändert =)(letztes oder -> und)

Aufgabe 4

(Auch hier ist die Problemstellung etwas unpräzise. Es wird keine Aussage gemacht, ob die "Autoknacker" und die "harmlosen Fußgängerkollisionen" (??) zusammen alle relevanten Erschütterungen ausmachen. Es wird ausnahmsweise mal vorausgesetzt.)

Es heißt doch, alle harmlosen Kollisionen, z.B. mit Spaziergängern –Philipendula 19:28, 11. Sep 2004 (UTC)

Die Wahrscheinlichkeiten, dass eine Erschütterung harmlos(h) bzw. durch einen Autoknacker (a) verursacht ist, betragen

- $w(h) = 0,8$
 - $w(a) = 0,2$
1. $w(b) = 0,9 * w(a) = 0,18 = 18\%$ **aller** Fälle (also Erschütterungen) bzw. $(0,9 * w(a)) / (0,9 * w(a) + 0,6 * w(h)) = 0,27... = 27,27\%$ aller Fälle, in denen Alarm ausgelöst wird (unklare Fragestellung!)
 2. $w(n) = (1-0,9) * w(a) + (1-0,6) * w(h) = 0,34 = 34\%$

Aufgabe 5

1. Besser, der Hausmeister stellt die Tonnen selber raus. Wozu ist er da?

Zum Tyrannisieren der Mieter und der Studies, die diese Aufgabe lösen müssen ;-) --[Philipendula](#) 19:30, 11. Sep 2004 (UTC)

- $w(L) = 0,3 * 0,07 = 0,021 = 2,1\%$
- $w(S) = 0,2 * 0,08 = 0,016 = 1,6\%$
- $w(F) = 0,5 * 0,05 = 0,025 = 2,5\%$ (--> am meisten verdächtig)

2. $w(Z) = w(L) + w(S) + w(F) = 6,2\%$

[Kategorie: Buch mit Druckversion](#)

Diese Seite wurde zuletzt am 30. Juni 2008 um 21:02 Uhr geändert.

Der Text ist unter der Lizenz „[Creative Commons Attribution/Share-Alike](#)“ verfügbar; zusätzliche Bedingungen können anwendbar sein. Einzelheiten sind in den [Nutzungsbedingungen](#) beschrieben.

[Datenschutz](#) [Über Wikibooks](#) [Impressum](#)

